

# Concept Space Alignment in Multilingual LLMs

Qiwei Peng and Anders Søgaard

University of Copenhagen

Denmark

{qipe, soegaard}@di.ku.dk

## Abstract

Multilingual large language models (LLMs) seem to generalize somewhat across languages. We hypothesize this is a result of implicit vector space alignment. Evaluating such alignment, we see that larger models exhibit *very* high-quality linear alignments between corresponding concepts in different languages. Our experiments show that multilingual LLMs suffer from two familiar weaknesses: generalization works best for languages with similar typology, and for abstract concepts. For some models, e.g., the Llama-2 family of models, prompt-based embeddings align better than word embeddings, but the projections are less linear – an observation that holds across almost all model families, indicating that some of the implicitly learned alignments are broken somewhat by prompt-based methods.

## 1 Introduction

Cross-lingual word embeddings are typically induced by supervised or unsupervised alignment of the word vector spaces of monolingual language models. Compression in multilingual models, i.e., parameter efficiency, can also drive *implicit* alignment (Devlin et al., 2019; Pires et al., 2019; Conneau et al., 2020), but until recently, the mappings could still be much improved by supervised or unsupervised alignment (Hu et al., 2021; Pan et al., 2021). Multilingual large language models (LLMs) are increasingly used for different tasks and demonstrate impressive ability in understanding different languages, but it is unclear whether this is a result of improved, implicit alignment, or of something else, e.g., linguistic overlap or semi-parallel subsets of training data.

LLMs have shown promising capability to comprehend *English* concepts (Liao et al., 2023; Xu et al., 2024). Our paper sets out to evaluate concept alignment in multilingual LLMs. We aim to investigate two things: First, is there a linear map-

en	failure, purchase, lizard, blink
fr	défaillance, emplette, lézard, ciller
ro	eșec, cumpărare, șopârlă, clipire
eu	porrot, erosketa, musker, betekara
fi	epäonnistuminen, osto, lisko, räpytys
ja	しくじり, 買いあげ, リザード, まぼたき
th	ความล้มเหลว, การซื้อ, ปอมข้าง, การกะพริบตา

Figure 1: Examples of four parallel WordNet concepts, aligned across 7 languages.

ping between corresponding concepts in different languages? Second, how does a learned linear mapping generalize to new concepts? We explore both questions by revisiting a set of techniques used in early work on bilingual dictionary induction (Kementchedjieva et al., 2018; Ruder et al., 2018; Søgaard et al., 2018; Kementchedjieva et al., 2019). We evaluate multilingual LLMs *as if* they were bilingual dictionary induction algorithms by doing nearest neighbor search – with cross-domain local scaling (Lample et al., 2018) – and evaluating their retrieval precision (precision@k). We first derive concept embedding in their standard way (last token or average). Since many of these models were *instruction fine-tuned*, we also compare *prompt-based embeddings* to standard techniques based on (low-level) word embeddings. We then compare their precision to retrieval rates after *explicit* concept space alignment. We perform analyses with and without leakage, across multiple languages, and across both abstract and physical concepts.

**Contributions** Our findings across experiments with 10 LLMs and six languages suggest that linear alignment can be induced in multilingual LLMs (if sufficiently big) to map concepts across different languages. Compared to vanilla embeddings, prompt-based concept embeddings exhibit significantly lower linearity, and the gaps between before and after alignment are larger for prompt-based

embeddings. This suggests that some of the implicitly learned concept alignments are broken by prompt-based methods. Prompt-based embeddings, which are now commonly used in different retrieval scenarios, seem to be less effective in extracting cross-lingually alignable embeddings, compared to vanilla embeddings. Results are generally good, but the old problem of generalization across typological distance (Singh et al., 2019) rears its ugly face again, with Basque, Finnish, Japanese and Thai exhibiting generally lower overall performance for both experimental set-ups. Furthermore, abstract concepts exhibit better alignment than physical concepts. We suspect that it is because abstract concepts are more frequent and occur in more diverse contexts.

## 2 Experiments

**Concepts** We collect English noun synsets from WordNet (Miller, 1995). For each synset, its first (most frequent) lemma name is used as the surface form of the corresponding concept. We use WordNet’s hierarchical structure to filter out top-level concepts (top-5 levels) to avoid too general concepts. WordNets in other languages, such as French WordNet (Sagot and Fišer, 2008), Basque WordNet (Gonzalez-Agirre et al., 2012), or Romanian WordNet (Dumitrescu et al., 2018), have similar structure and were all aligned in the Open Multilingual WordNet project (OMW) (Bond et al., 2016). To produce a repository of parallel semantic concepts, we collect synsets with shared ID across different WordNets, after removing duplicate concepts. In total, we obtain 4,397 parallel concepts across 7 different languages (English, French, Romanian, Basque, Finnish (Lindén and Niemi, 2014), Japanese (Bond et al., 2009) and Thai (Thoongsup et al., 2009)); had we included more languages, the number of parallel concepts would have been prohibitively small. The 4,397 concepts were divided into abstract (e.g., happiness) and physical (e.g., vehicle) concepts. See Table 1 for data characteristics, and Figure 1 for examples of parallel concepts.

	Abstract	Physical	Total
Train	1500	1500	3000
Test	419	978	1397
Total	1919	2478	4397

Table 1: The statistics of the parallel concept dataset. We use 1000, 2000, or 3000 concepts for training.

To create a **seed dictionary** (training data) for supervised alignment, we randomly sample 3,000 parallel concepts,<sup>1</sup> including 1,500 abstract concepts and 1,500 physical concepts. The 3,000 concepts are used to induce the linear mapping.

**LLMs** We experiment with four different LLM families with varying sizes: Llama2 (7B, 13B, 70B) (Touvron et al., 2023), mT0 (1.2B, 3.7B, 13B), BLOOMZ (1B7, 3B, 7B1) (Muennighoff et al., 2022), and Aya101 (13B) (Üstün et al., 2024). We use *two different* concept space extraction methods (vanilla and prompt-based). The vanilla method simply uses the last token representation as the concept embedding for decoder-only models (Llama2 and BLOOMZ); and the average embedding of the last hidden layer of the encoder as the concept embedding for encoder-decoder models (mT0 and Aya101)<sup>2</sup>. The prompt-based extraction method exploits the fact that all these models were instruction-tuned. The template we use for prompt-based extraction is adapted from Li and Li (2023) and shown as follows:

Summarize concept [text] in one [lang] word:

where [text] and [lang] will be replaced by the corresponding concept (in the source language) and the language name (in adjectival form), e.g., "summarize concept "動物" in one Japanese word" for the concept *animal*<sup>3</sup>. The prompt-based concept embedding is that of the last hidden state.

**Alignment and Retrieval** We rely on Procrustes Analysis (Schönemann, 1966), a form of statistical shape analysis, to discover good linear transformations (e.g., translation, rotation, and scaling) between concept spaces in different languages. Suppose  $X$  and  $Y$  are two matrices of size  $n \times d$  ( $n$  is the seed dictionary size,  $d$  is the embedding size) such that the  $i$ th row in  $X$  is an embedding of concept  $c_i$  in one language, and the  $i$ th row in  $Y$  is  $c_i$ ’s embedding in the other language. The linear transformation is derived through singular value decomposition (SVD) of  $YX^T$ :

$$W^* = \arg \min_{W \in O_d(\mathbb{R})} \|WX - Y\|_F = UV^T \quad (1)$$

where  $U\Sigma V^T = \text{SVD}(YX^T)$ . With  $W^*$ , we transform source language concept embeddings  $X$  into

<sup>1</sup>See Appendix for results with 1,000 or 2,000 concepts.

<sup>2</sup>This is decided by preliminary experiment results.

<sup>3</sup>The [text] of a concept can be made up of multiple words.

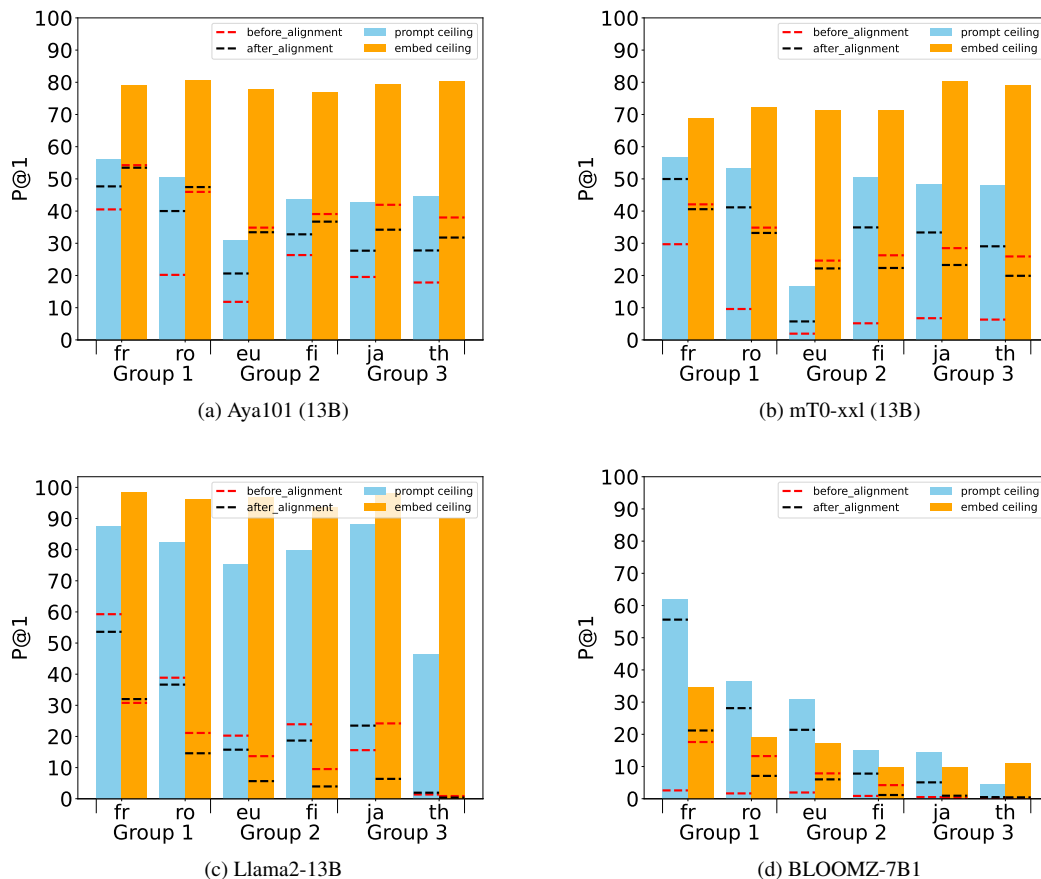


Figure 2: Performance (P@1) of different LLMs on the concept alignment evaluation when using a seed dictionary of 3,000 concepts. X-axis: Languages, we further divide these languages into three groups, where **Group 1** is Indo-European, **Group 2** includes languages that are not Indo-European but still in Latin script, while **Group 3** refers to languages that are not Indo-European and not in Latin script. Y-axis: We report Precision@1.

the *English* (target) vector space. We then perform cross-domain local scaling (CSLS) to retrieve the most similar concepts.<sup>4</sup> We use precision@k (P@k) as our performance metric.

**Main Results** We present the main results<sup>5</sup> in Figure 2. For each model, we report three results: 1) the *upper bound* (leaky) on performance for supervised linear alignment, using the train seed and the test seed for inducing the dictionary (orange/blue bar), which we refer to as the *ceiling* and reveals to what extent there exists a linear mapping; 2) *before-align* performance (red dashed line), retrieval bilingual concept pairs directly from the raw LLM (vanilla word or prompt) embeddings; 3) *after-align* performance (black dashed line), which

<sup>4</sup>We also ran experiments with vanilla nearest neighbor search as our retrieval method, but CSLS outperforms nearest neighbor search by some margin. So, we report results with CSLS.

<sup>5</sup>Full results with all model sizes, training sizes, and different  $k$ -values for P@k are presented in the Appendix.

is the performance of non-leaky, supervised mapping (using 3,000 concepts as the seed dictionary) into the English vector space, with CSLS as our retrieval method. Orange bars indicate vanilla word embedding strategy (last-token, or average embedding), while blue bars refer to results for prompt-based embedding.

All multilingual LLMs (except BLOOMZ) can induce good concept alignments, as indicated by the upper bound performance. In general, within the same model family, a larger model size leads to better alignment. The ceiling is highest for vanilla word embeddings in Llama2-13B, indicating near-isomorphisms between monolingual concept spaces at this level. The prompt-based embeddings are less linear, indicating that partial isomorphisms induced prior to prompting are corrupted. For after-align performance, we generally see the highest performance for Indo-European languages (Group 1) and the lowest for non-Indo-European languages with non-Latin scripts (Group 3). Sim-

ilarly, a larger model size and larger seed dictionary generally improve the concept alignment. On Group 2 and 3, mT0 and Aya101 show better before-align performance compared to other models. In some cases, results are extremely good. Llama2-13B with prompt-based embeddings exhibits a P@1 score of 59.27% before alignment for French, for example. This means that the model has induced perfect alignment of 3/5 concepts in the absence of any explicit supervision. It is interesting to see the gap between the red and black dashed lines. The size of this gap indicates how much of the (alignable part of the) concept space was *not* aligned, with given seed dictionary. For vanilla word embeddings, the gaps are relatively small, but for prompt-based embeddings the gaps tend to be much larger, again indicating that prompting somewhat breaks the implicitly learned concept alignment.

**Abstract vs. Physical** We analyze performance differences across abstract and physical concepts. To make a fair comparison, we randomly down-sample<sup>6</sup> physical concepts and compare retrieval performance across the two classes. In this section, we report P@1 with models that have comparable model sizes (7B/13B) in each family; results for the other models can be found in the Appendix. As shown in Table 2, all models generally have better alignment performance on abstract concepts compared to physical concepts.

		fr	ro	eu	fi	ja	th
Llama2-13B	Abstract	<b>63.48</b>	<b>46.06</b>	<b>17.42</b>	<b>21.00</b>	<b>26.01</b>	<b>2.15</b>
	Physical	50.12	33.41	14.08	18.62	23.39	1.91
BLOOMZ-7B1	Abstract	<b>64.92</b>	<b>33.41</b>	<b>27.92</b>	<b>10.74</b>	<b>10.26</b>	<b>1.19</b>
	Physical	52.51	27.45	18.62	6.44	4.53	0.00
mT0-xxl (13B)	Abstract	<b>59.90</b>	<b>49.88</b>	<b>7.88</b>	<b>38.90</b>	<b>38.42</b>	<b>34.37</b>
	Physical	46.78	41.29	5.73	37.23	36.28	28.64
Aya101 (13B)	Abstract	<b>58.47</b>	<b>52.27</b>	<b>27.68</b>	<b>40.81</b>	<b>36.28</b>	<b>32.70</b>
	Physical	44.63	36.75	18.38	30.79	26.73	29.12

Table 2: The results (P@1) for abstract and physical concepts. We report after-align results for prompt-based embedding and comparable sizes (13B/7B) of each model family.

What explains this very consistent finding? One hypothesis would be that physical nouns are more ambiguous, since they often source metaphor and metonymy. However, our words for abstract concepts have more senses in WordNet (2.94) than our words for physical concepts (1.96); see Table

<sup>6</sup>See Appendix for numbers without down-sampling.

3. Instead, we found another, simpler explanation. Frequency statistics (obtained from the English Wikipedia dump of 2023-04-13) relevant that the abstract concept words are considerably more frequent than the physical concept words, which makes sense, as abstract concepts apply very generally across contexts and domains.

	Abstract	Physical
avg # of senses	2.94	1.96
median # of senses	2	1
avg # of counts	103,934	28,762
median # of count	12,787	5,122

Table 3: Number of senses and frequency of words.

### 3 Discussion and Related Work

**Related Work** The idea that distributional representations facilitate cross-lingual alignment goes back to explicit semantic analysis (Gabrilovich and Markovitch, 2007), but the idea of training multilingual, neural language models also has a long history. Such models have traditionally used explicit alignment objectives, e.g., either from word alignments, bilingual dictionary seeds (Lample et al., 2018; Li et al., 2024), or by training on mixed corpora constructed using such resources (Gouws and Søgaard, 2015; Workshop et al., 2022; Chai et al., 2024). Cross-lingual generalization has been studied in different NLP tasks, including question answering (Artetxe et al., 2020), commonsense reasoning (Ponti et al., 2020; Lin et al., 2022), code generation (Peng et al., 2024), and knowledge transfer and consistency (Xu et al., 2023; Qi et al., 2023). Cross-lingual word alignment also has a long history by examining bilingual lexicon induction (Xing et al., 2015; Søgaard et al., 2018; Li et al., 2023). For concept understanding specifically, previous works have examined concept understanding in LLMs by definition matching (Xu et al., 2024), hypernym/hyponym detection (Liao et al., 2023; Shani et al., 2023), and relation discovery (Gu et al., 2023). However, they are limited to the English language only.

**Linear Alignment** We saw that concepts are represented in similar ways across languages in multilingual LLMs, as shown in the upper bound. This indicates structural similarities and facilitates cross-lingual transfer. Prompt-based embeddings exhibit significantly lower linearity compared to word embeddings, and the gaps between before and after

alignment are larger for prompt-based embeddings. Both things suggest that some of the implicitly learned concept alignment is broken by the prompt-based method. On the other hand, prompt-based embeddings demonstrate larger improvements with explicit post-hoc alignment while supervised alignment struggles to improve on vanilla word embeddings.

**Difference in Languages** The degree of isomorphism to English is similar across languages, as indicated by the upper bounds on performance. All concept spaces are (almost) equally alignable. However, the induced maps generalize much better across typologically related (Indo-European) languages: French and Romanian. Generalization is considerably poorer for the other two groups.

**Types of Concepts** Though previous works show that physical concepts do better than abstract ones in bilingual dictionary induction (Kementchedjheva et al., 2019), as well as in related tasks such as hypernym detection (Liao et al., 2023), we show that abstract concepts tend to align better across different languages, as shown in Table 2. This, however, was explained by a spurious correlation with frequency. It would be interesting to control for frequency in future error analysis.

## 4 Conclusion

We evaluated concept alignment on multilingual LLMs by revisiting the traditional bilingual dictionary induction task, but with semantic concepts rather than words. Our experiments show that multilingual LLMs exhibit high-quality, linear concept alignment across different languages. However, the ability of supervised maps to generalize varied across different models, languages, and ways of obtaining embeddings.

## Limitations

Because of the small overlap between multilingual WordNets, we only include six (6) test languages. While this is too small a set of languages to draw universally applicable conclusions. Fortunately, the set includes both Indo-European and non-Indo-European languages, as well as both Latin script and non-Latin script languages. We also limited ourselves to studying nouns; for how linear alignment generalizes to other parts of speech, see Kementchedjheva et al. (2018) and Hartmann and Søgaard (2018).

## Ethical Considerations

We do not anticipate any risks in the work. In this study, our use of existing artifacts is consistent with their intended purposes. Semantic concepts are collected from previously published and publicly available resources (WordNets). Aya101<sup>7</sup>, BLOOMZ, and mT0 models have Apache-2.0 License<sup>8</sup>. Llama2 models are under the LLAMA 2 Community License<sup>9</sup>.

## Acknowledgement

We would like to thank all anonymous reviewers for their insightful comments and feedback. This work was supported by DisAI - Improving scientific excellence and creativity in combating disinformation with artificial intelligence and language technologies, a project funded by European Union under the Horizon Europe, GA No. 101079164.

## References

- Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. 2020. On the cross-lingual transferability of monolingual representations. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4623–4637.
- Francis Bond, Hitoshi Isahara, Sanae Fujita, Kiyotaka Uchimoto, Takayuki Kuribayashi, and Kyoko Kan-zaki. 2009. Enhancing the japanese wordnet. In *Proceedings of the 7th Workshop on Asian Language Resources (ALR7)*, pages 1–8.
- Francis Bond, Piek Vossen, John Philip McCrae, and Christiane Fellbaum. 2016. Cili: the collaborative interlingual index. In *Proceedings of the 8th Global WordNet Conference (GWC)*, pages 50–57.
- Linzhen Chai, Jian Yang, Tao Sun, Hongcheng Guo, Jiaheng Liu, Bing Wang, Xiannian Liang, Jiaqi Bai, Tongliang Li, Qiyao Peng, et al. 2024. xcot: Cross-lingual instruction tuning for cross-lingual chain-of-thought reasoning. *arXiv preprint arXiv:2401.07037*.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. [Unsupervised cross-lingual representation learning at scale](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.
- 
- <sup>7</sup><https://huggingface.co/CohereForAI/aya-101#model-summary>
- <sup>8</sup><https://github.com/bigscience-workshop/xmft/blob/master/README.md>
- <sup>9</sup><https://ai.meta.com/llama/license/>

- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). 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 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Stefan Daniel Dumitrescu, Andrei Marius Avram, Luciana Morogan, and Stefan-Adrian Toma. 2018. Rowordnet—a python api for the romanian wordnet. In *2018 10th International Conference on Electronics, Computers and Artificial Intelligence (ECAI)*, pages 1–6. IEEE.
- Evgeniy Gabrilovich and Shaul Markovitch. 2007. Computing semantic relatedness using wikipedia-based explicit semantic analysis. In *Proceedings of the 20th International Joint Conference on Artificial Intelligence*, volume 7, pages 1606–1611.
- Aitor Gonzalez-Agirre, Egoitz Laparra, German Rigau, et al. 2012. Multilingual central repository version 3.0. In *LREC*, pages 2525–2529.
- Stephan Gouws and Anders Søgaard. 2015. [Simple task-specific bilingual word embeddings](#). In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1386–1390, Denver, Colorado. Association for Computational Linguistics.
- Yuling Gu, Bhavana Dalvi, and Peter Clark. 2023. Do language models have coherent mental models of everyday things? In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1892–1913.
- Mareike Hartmann and Anders Søgaard. 2018. [Limitations of cross-lingual learning from image search](#). In *Proceedings of the Third Workshop on Representation Learning for NLP*, pages 159–163, Melbourne, Australia. Association for Computational Linguistics.
- Junjie Hu, Melvin Johnson, Orhan Firat, Aditya Siddhant, and Graham Neubig. 2021. [Explicit alignment objectives for multilingual bidirectional encoders](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3633–3643, Online. Association for Computational Linguistics.
- Yova Kementchedjhieva, Mareike Hartmann, and Anders Søgaard. 2019. [Lost in evaluation: Misleading benchmarks for bilingual dictionary induction](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3336–3341, Hong Kong, China. Association for Computational Linguistics.
- Yova Kementchedjhieva, Sebastian Ruder, Ryan Cotterell, and Anders Søgaard. 2018. [Generalizing Procrustes analysis for better bilingual dictionary induction](#). In *Proceedings of the 22nd Conference on Computational Natural Language Learning*, pages 211–220, Brussels, Belgium. Association for Computational Linguistics.
- Guillaume Lample, Alexis Conneau, Marc’Aurelio Ranzato, Ludovic Denoyer, and Hervé Jégou. 2018. [Word translation without parallel data](#). In *International Conference on Learning Representations*.
- Jiahuan Li, Shujian Huang, Xinyu Dai, and Jiajun Chen. 2024. Prealign: Boosting cross-lingual transfer by early establishment of multilingual alignment. *arXiv preprint arXiv:2407.16222*.
- Xianming Li and Jing Li. 2023. Angle-optimized text embeddings. *arXiv preprint arXiv:2309.12871*.
- Yaoyiran Li, Anna Korhonen, and Ivan Vulić. 2023. [On bilingual lexicon induction with large language models](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9577–9599, Singapore. Association for Computational Linguistics.
- Jiayi Liao, Xu Chen, and Lun Du. 2023. [Concept understanding in large language models: An empirical study](#). In *Tiny Papers @ ICLR*.
- Xi Victoria Lin, Todor Mihaylov, Mikel Artetxe, Tianlu Wang, Shuohui Chen, Daniel Simig, Myle Ott, Naman Goyal, Shruti Bhosale, Jingfei Du, et al. 2022. Few-shot learning with multilingual generative language models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9019–9052.
- Krister Lindén and Jyrki Niemi. 2014. Is it possible to create a very large wordnet in 100 days? an evaluation. *Language resources and evaluation*, 48:191–201.
- George A Miller. 1995. Wordnet: a lexical database for english. *Communications of the ACM*, 38(11):39–41.
- Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hailey Schoelkopf, et al. 2022. Crosslingual generalization through multitask finetuning. *arXiv preprint arXiv:2211.01786*.
- Lin Pan, Chung-Wei Hang, Haode Qi, Abhishek Shah, Saloni Potdar, and Mo Yu. 2021. [Multilingual BERT post-pretraining alignment](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 210–219, Online. Association for Computational Linguistics.
- Qiwei Peng, Yekun Chai, and Xuhong Li. 2024. Humaneval-xl: A multilingual code generation

- benchmark for cross-lingual natural language generalization. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 8383–8394.
- Telmo Pires, Eva Schlinger, and Dan Garrette. 2019. [How multilingual is multilingual BERT?](#) In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4996–5001, Florence, Italy. Association for Computational Linguistics.
- Edoardo Maria Ponti, Goran Glavaš, Olga Majewska, Qianchu Liu, Ivan Vulić, and Anna Korhonen. 2020. Xcopa: A multilingual dataset for causal commonsense reasoning. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2362–2376.
- Jirui Qi, Raquel Fernández, and Arianna Bisazza. 2023. [Cross-lingual consistency of factual knowledge in multilingual language models.](#) In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 10650–10666, Singapore. Association for Computational Linguistics.
- Sebastian Ruder, Ryan Cotterell, Yova Kementchedzhieva, and Anders Søgaard. 2018. [A discriminative latent-variable model for bilingual lexicon induction.](#) In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 458–468, Brussels, Belgium. Association for Computational Linguistics.
- Benoît Sagot and Darja Fišer. 2008. Building a free french wordnet from multilingual resources. In *OntoLex*.
- Peter H Schönemann. 1966. A generalized solution of the orthogonal procrustes problem. *Psychometrika*, 31(1):1–10.
- Chen Shani, Jilles Vreeken, and Dafna Shahaf. 2023. Towards concept-aware large language models. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 13158–13170.
- Jasdeep Singh, Bryan McCann, Richard Socher, and Caiming Xiong. 2019. [BERT is not an interlingua and the bias of tokenization.](#) In *Proceedings of the 2nd Workshop on Deep Learning Approaches for Low-Resource NLP (DeepLo 2019)*, pages 47–55, Hong Kong, China. Association for Computational Linguistics.
- Anders Søgaard, Sebastian Ruder, and Ivan Vulić. 2018. [On the limitations of unsupervised bilingual dictionary induction.](#) In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 778–788, Melbourne, Australia. Association for Computational Linguistics.
- Sareewan Thoongsup, Thatsanee Charoenporn, Kergit Robkop, Tan Sinthurahat, Chumpol Mokrat, Virach Sornlertlamvanich, and Hitoshi Isahara. 2009. Thai wordnet construction. In *Proceedings of the 7th Workshop on Asian Language Resources (ALR7)*, pages 139–144.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shriti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Ahmet Üstün, Viraat Aryabumi, Zheng-Xin Yong, Wei-Yin Ko, Daniel D’souza, Gbemileke Onilude, Neel Bhandari, Shivalika Singh, Hui-Lee Ooi, Amr Kayid, et al. 2024. Aya model: An instruction finetuned open-access multilingual language model. *arXiv preprint arXiv:2402.07827*.
- BigScience Workshop, Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luciani, François Yvon, et al. 2022. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*.
- Chao Xing, Dong Wang, Chao Liu, and Yiye Lin. 2015. [Normalized word embedding and orthogonal transform for bilingual word translation.](#) In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1006–1011, Denver, Colorado. Association for Computational Linguistics.
- Ningyu Xu, Qi Zhang, Menghan Zhang, Peng Qian, and Xuanjing Huang. 2024. On the tip of the tongue: Analyzing conceptual representation in large language models with reverse-dictionary probe. *arXiv preprint arXiv:2402.14404*.
- Shaoyang Xu, Junzhuo Li, and Deyi Xiong. 2023. [Language representation projection: Can we transfer factual knowledge across languages in multilingual language models?](#) In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 3692–3702, Singapore. Association for Computational Linguistics.

## A Language Resource and Shared Vocabulary

We report the estimated resource level for the seven languages we experimented in this work. The number, which has been widely used to indicate resource availability, is taken from the CC100 XL corpus Lin et al. (2022).

One reason of the cross-lingual alignment might be that among the chosen concepts some had the same word form across languages with Latin scripts. To investigate this, we additionally calculate the ratio of identical word forms (compared to English concepts) for languages with Latin scripts:

Lang	Tokens (M)
English	803,527
French	77,420
Japanese	66,054
Romanian	24,176
Finnish	16,804
Thai	10,842
Basque	105

Table 4: The statistics of language resource level for 7 languages used in this work.

- 653 out of 4397 for French (roughly 15%)
- 449 out of 4397 for Romanian (roughly 10%)
- 243 out of 4397 for Basque (roughly 5%)
- 172 out of 4397 for Finnish (roughly 4%)

This ratio remains similar when counted only on the test split. We can observe from above that the ratio is quite limited.

## B Full Experimental Results

In the Appendix, we report our full experimental results across different models with varying model sizes, seed dictionary sizes, different k-values for P@K, in following figure and tables (Figure 3 and Table 4-23). These results provide full scope of our analysis, allowing for an in-depth comparison of model performances.

For different models, we use their HuggingFace PyTorch implementation<sup>10</sup>. For Procrustes Analysis, we utilize the MUSE<sup>11</sup> package. All experiments are run on a single NVIDIA A100 GPU.

<sup>10</sup><https://huggingface.co/bigscience/bloomz-1b7,3b,7b1>,  
<https://huggingface.co/meta-llama/Llama-2-7,13,70b-chat-hf>,  
<https://huggingface.co/bigscience/mt0-large,xl,xxl>,  
<https://huggingface.co/CoHereForAI/aya-101>

<sup>11</sup><https://github.com/facebookresearch/MUSE>



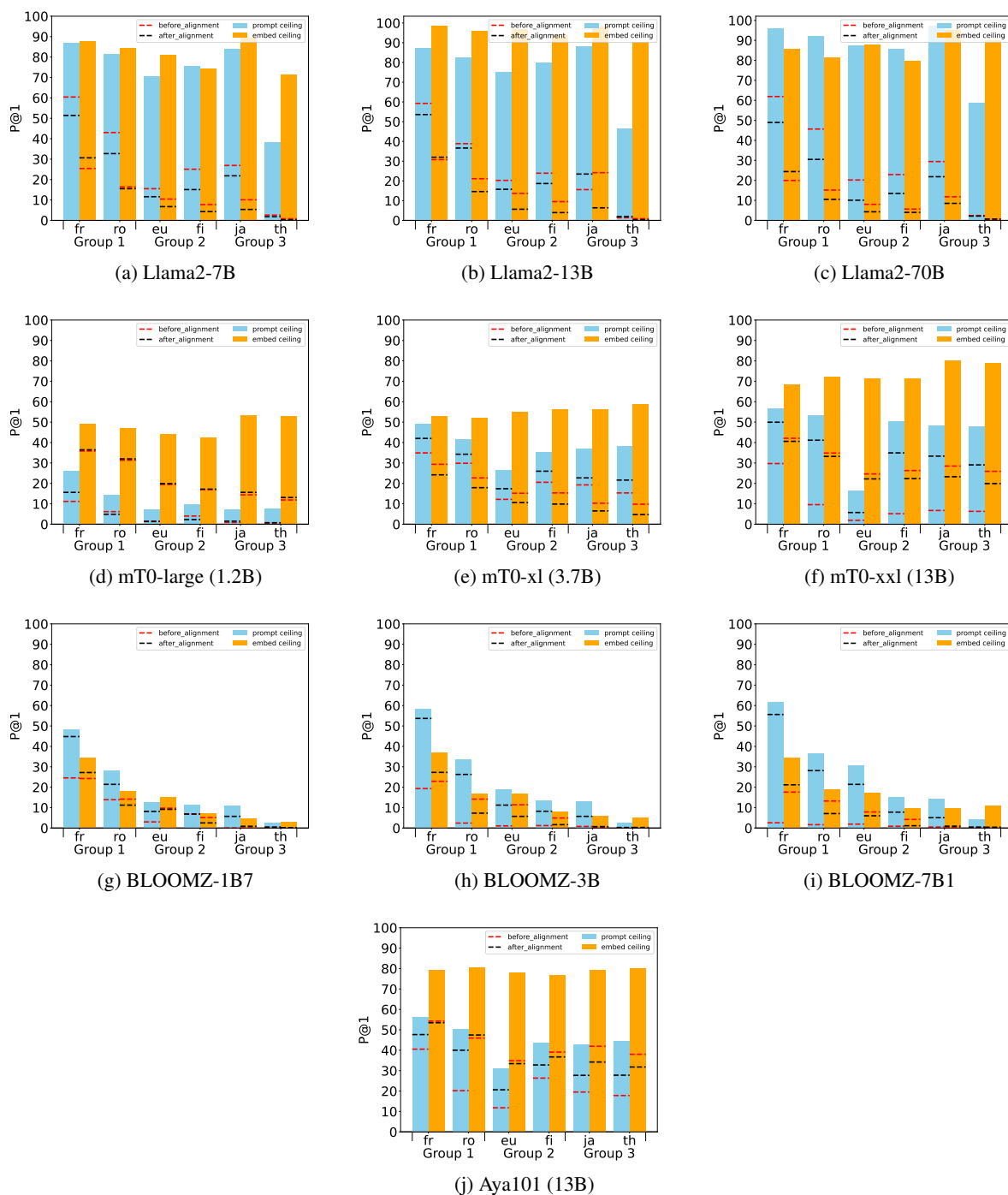


Figure 3: Performance (P@1) of different LLMs on the concept alignment evaluation when using a seed dictionary of 3000 pairs. X-axis: Languages, we further divide these languages into three groups, where **Group 1** is Indo-European, **Group 2** includes languages that are not Indo-European but still in Latin script, while **Group 3** refers to languages that are not Indo-European and not in Latin script. Y-axis: We report Precision@1.













EN (P@K: 1000/2000/3000)							
Language	P@K	Before-Alignment	Ceiling Performance	Eval - Both	Eval - Abstract	Eval-Physical (all)	Eval - Physical (downsample)
fr	p@1	19.97	85.90	13.17/20.04/24.41	14.08/22.20/25.54	12.78/19.12/23.93	11.22/17.66/23.87
	p@5	23.62	89.33	27.06/34.79/38.30	27.92/38.19/42.24	26.69/33.33/36.61	26.97/33.41/36.75
	p@10	26.63	90.41	32.07/39.58/42.59	35.80/44.39/46.78	30.47/37.53/40.80	31.03/36.99/42.48
	p@30	30.57	92.13	39.44/45.88/48.32	44.39/51.55/52.98	37.32/43.46/46.32	38.42/44.39/48.45
ro	p@1	15.18	81.39	4.51/7.52/10.52	3.34/6.44/10.26	5.01/7.98/10.63	3.82/7.64/10.98
	p@5	17.47	85.25	14.39/18.90/21.55	15.04/19.09/21.72	14.11/18.81/21.47	15.04/19.81/23.15
	p@10	18.97	87.33	18.90/24.27/26.63	20.29/25.78/28.88	18.30/23.62/25.66	19.09/24.34/27.21
	p@30	22.55	90.12	26.77/31.21/34.72	29.83/33.89/39.38	25.46/30.06/32.72	26.25/31.98/34.61
eu	p@1	8.02	87.97	1.36/2.43/4.29	1.91/2.86/3.58	1.12/2.25/4.60	0.72/2.39/4.77
	p@5	9.81	91.62	7.37/11.02/13.17	8.35/11.93/15.04	6.95/10.63/12.37	6.92/11.93/12.89
	p@10	11.60	92.63	10.16/14.53/16.39	11.69/16.95/18.62	9.51/13.50/15.44	9.55/15.04/15.99
	p@30	14.46	94.42	16.46/20.69/22.41	19.57/23.87/25.54	15.13/19.33/21.06	14.32/20.76/22.20
fi	p@1	5.65	79.96	1.57/2.51/4.01	2.63/2.86/4.53	1.12/2.35/3.78	0.72/0.95/2.63
	p@5	8.45	85.18	7.52/10.74/13.17	10.02/12.89/15.75	6.44/9.82/12.07	5.97/10.26/12.41
	p@10	10.02	86.61	12.31/15.25/18.04	15.99/17.90/20.76	10.74/14.11/16.87	10.50/15.27/18.62
	p@30	14.53	88.83	19.76/22.69/24.84	24.58/25.78/28.88	17.69/21.37/23.11	17.90/22.67/23.87
ja	p@1	11.81	95.13	3.15/6.16/8.52	3.82/7.64/10.02	2.86/5.52/7.87	3.58/7.16/9.31
	p@5	22.05	96.99	12.46/19.33/23.41	12.65/20.29/25.54	12.37/18.92/22.49	13.37/20.53/23.39
	p@10	25.70	97.28	18.25/25.48/29.63	19.09/26.97/31.98	17.89/24.85/28.63	20.29/26.25/29.36
	p@30	31.78	98.35	27.63/35.43/39.58	30.31/37.23/40.33	26.48/34.66/39.26	28.16/36.28/40.81
th	p@1	0.64	89.91	0.21/0.43/0.64	0.24/0.24/0.48	0.20/0.51/0.72	0.24/1.19/1.91
	p@5	1.57	92.98	2.72/4.58/6.23	2.39/4.77/6.92	2.86/4.50/5.93	3.58/5.25/6.68
	p@10	2.58	93.49	4.72/6.30/8.95	5.49/6.68/10.74	4.40/6.13/8.18	4.53/6.44/8.35
	p@30	4.37	94.56	10.31/13.53/15.39	11.93/15.99/17.18	9.61/12.47/14.62	9.55/12.89/14.80

Table 23: Full results for Llama2-70B with last-token embedding

EN (P@K: 1000/2000/3000)							
Language	P@K	Before-Alignment	Ceiling Performance	Eval - Both	Eval - Abstract	Eval-Physical (all)	Eval - Physical (downsample)
fr	p@1	61.92	96.06	34.00/44.31/48.96	42.00/53.22/56.56	30.57/40.49/45.71	30.31/41.29/45.35
	p@5	72.15	99.36	66.43/69.86/71.44	74.22/77.80/79.47	63.09/66.46/68.00	65.63/68.26/69.69
	p@10	74.66	99.50	71.73/74.16/75.02	78.52/81.38/81.38	68.81/71.06/72.29	71.84/73.03/74.22
	p@30	78.53	99.93	77.17/78.95/79.96	83.05/84.25/84.49	74.64/76.69/78.02	77.57/79.00/80.19
ro	p@1	45.67	92.20	16.96/26.56/30.49	22.43/36.52/37.95	14.62/22.29/27.30	13.84/21.48/26.01
	p@5	55.19	97.71	47.96/53.97/55.33	55.13/60.62/61.81	44.89/51.12/52.56	45.35/51.07/52.98
	p@10	57.84	98.57	54.76/58.91/60.63	61.58/66.11/68.74	51.84/55.83/57.16	52.74/55.85/57.04
	p@30	61.27	99.50	63.42/65.43/66.86	73.27/73.51/74.46	59.20/61.96/63.60	59.90/63.25/65.16
eu	p@1	20.19	87.54	3.65/6.87/10.09	7.16/10.50/13.84	2.15/5.32/8.49	2.86/6.92/9.79
	p@5	27.49	95.28	19.76/26.34/29.13	23.63/31.50/33.41	18.10/24.13/27.30	19.57/25.06/27.45
	p@10	29.78	96.85	24.34/30.06/33.21	26.73/34.61/36.52	23.31/28.12/31.80	25.54/29.12/32.94
	p@30	31.93	98.14	30.49/36.08/38.58	33.65/41.77/42.96	29.14/33.64/36.71	30.55/35.08/38.19
fi	p@1	22.91	85.61	6.01/10.52/13.46	9.07/15.99/21.00	4.70/8.18/10.22	4.53/8.59/10.26
	p@5	29.21	92.70	27.77/33.86/36.36	35.32/42.00/43.91	24.54/30.37/33.13	26.01/31.26/34.13
	p@10	31.85	94.56	33.50/39.16/41.59	41.77/47.26/49.40	29.96/35.69/38.24	32.22/37.23/39.62
	p@30	34.72	96.49	42.09/45.74/47.17	49.40/53.70/54.65	38.96/42.33/43.97	42.24/44.87/45.82
ja	p@1	29.42	97.21	12.67/18.18/21.83	14.56/21.00/23.87	11.86/16.97/20.96	13.13/17.90/21.72
	p@5	44.74	99.36	44.02/51.83/55.69	52.51/58.71/62.77	40.39/48.88/52.66	44.87/53.94/57.28
	p@10	51.11	99.43	53.33/60.77/63.85	61.34/68.02/71.60	49.90/57.67/60.53	53.70/61.34/65.16
	p@30	60.49	99.86	65.00/70.87/73.30	72.55/78.28/80.19	61.76/67.69/70.35	64.68/69.69/72.79
th	p@1	2.43	58.98	1.07/1.43/2.15	0.48/1.19/2.63	1.33/1.53/1.94	1.67/1.67/1.91
	p@5	5.44	74.80	7.16/9.88/11.52	8.59/11.69/14.56	6.54/9.10/10.22	7.16/8.59/9.31
	p@10	7.66	79.24	11.81/16.25/17.32	13.60/20.05/20.76	11.04/14.62/15.85	11.22/14.80/15.51
	p@30	12.03	86.54	20.69/25.05/27.27	26.25/30.31/30.55	18.30/22.80/25.87	16.95/21.72/24.34

Table 24: Full results for Llama2-70B with prompt-based embedding