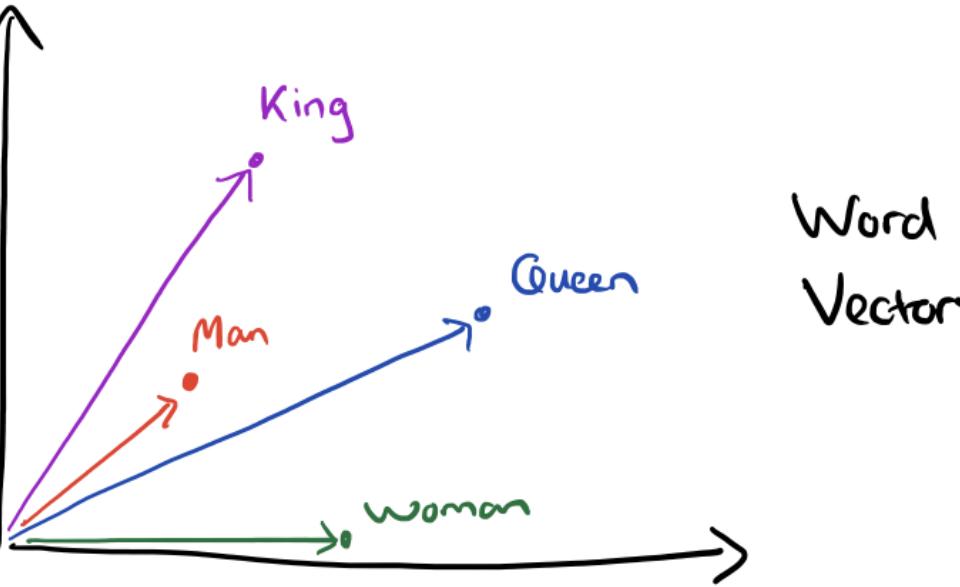
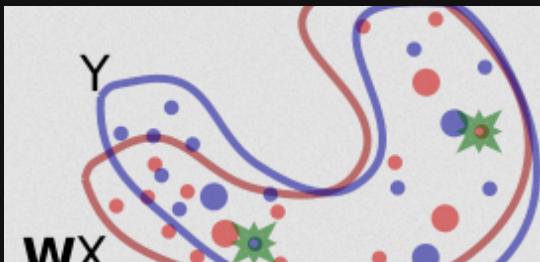


{Multilingual Intertext Detection}

in Ancient Languages

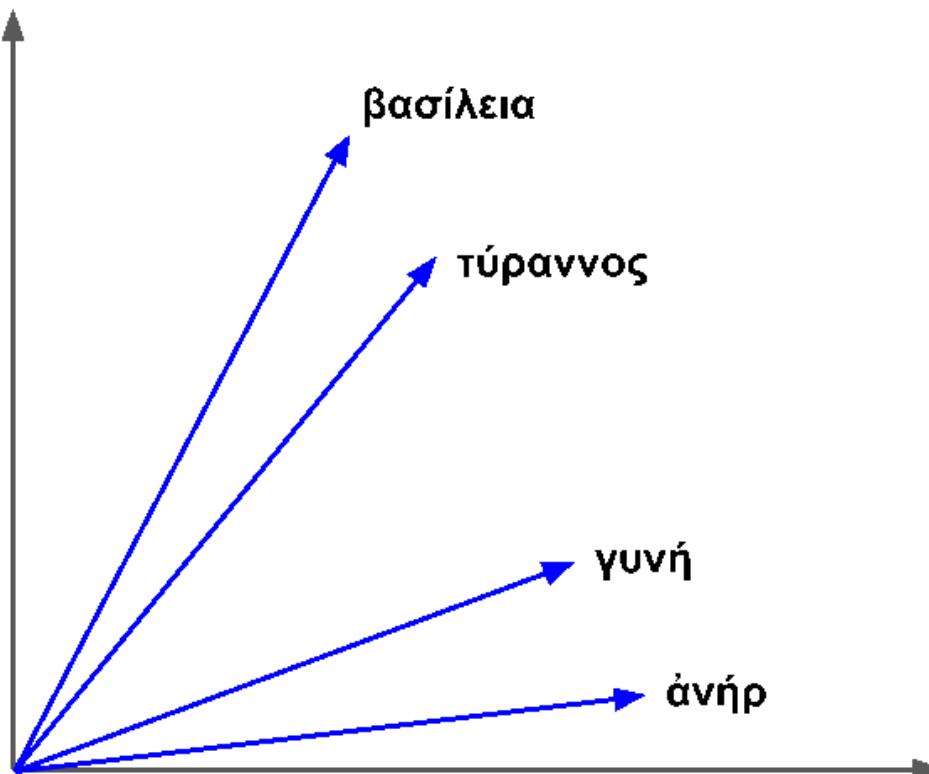


Basic concepts



static word vectors, contextual vectors, ...

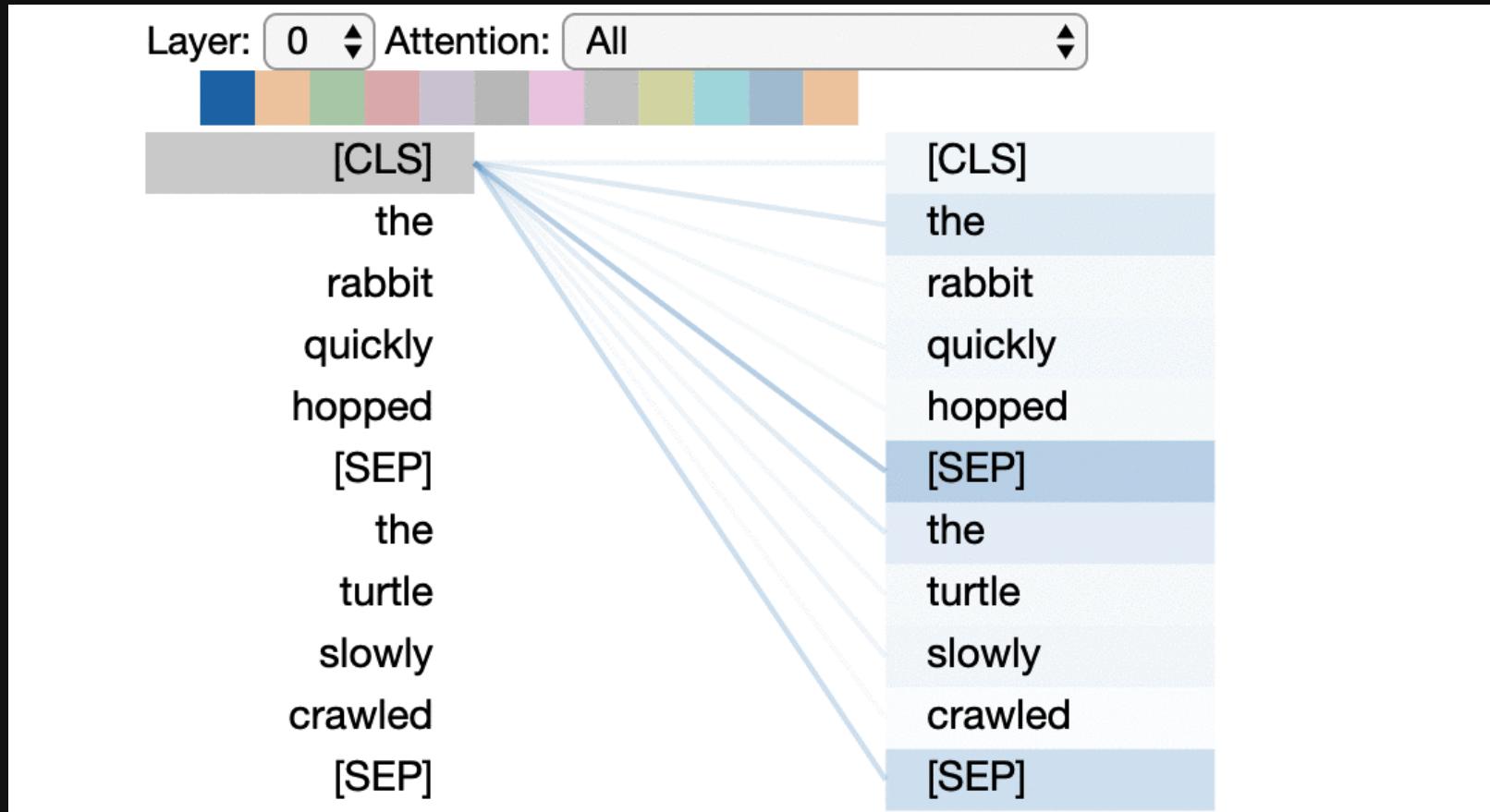
Word vectors



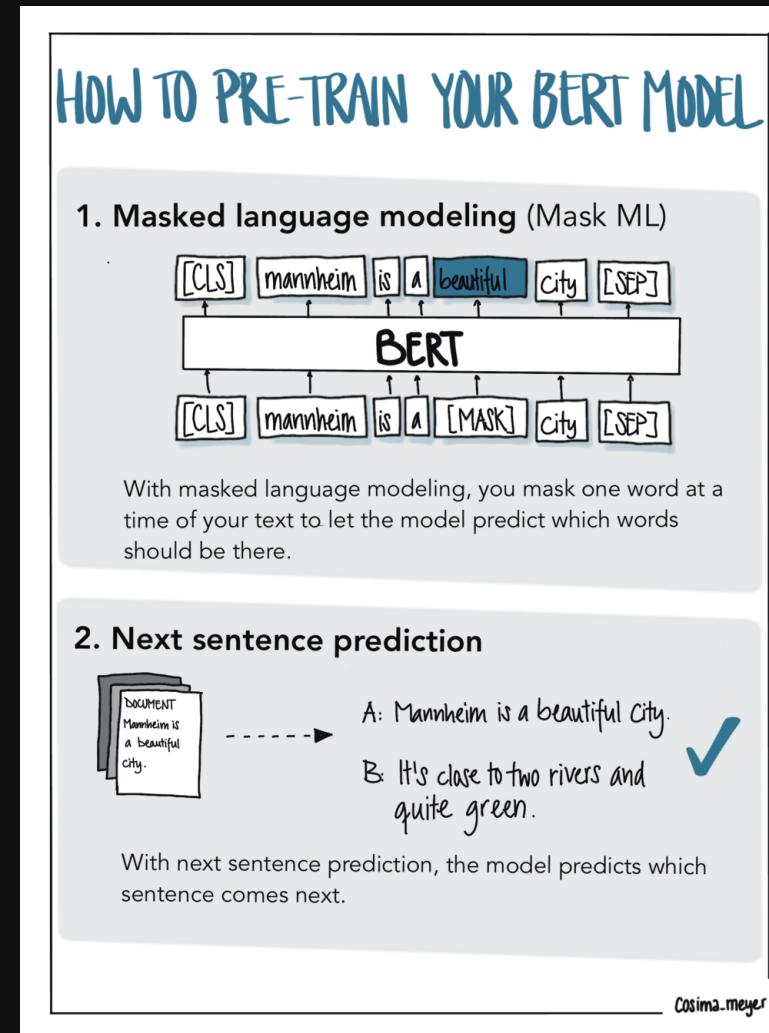
Example on a medium corpus (Balzac)

<https://projector.tensorflow.org/>

What's a large language model (LLM) ?



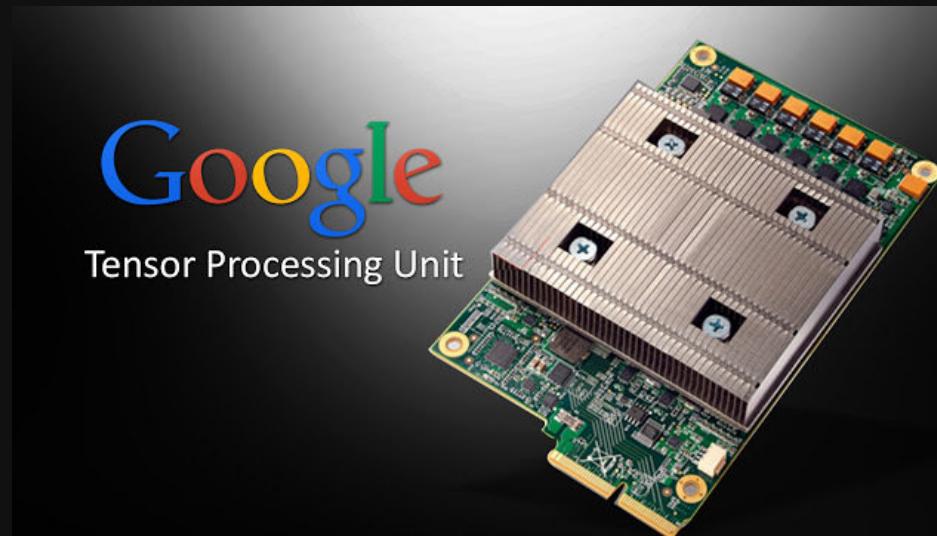
Training



<https://www.mzes.uni-mannheim.de/socialsciencedatalab/article/bert-explainable-ai/#bert>

Problems ?

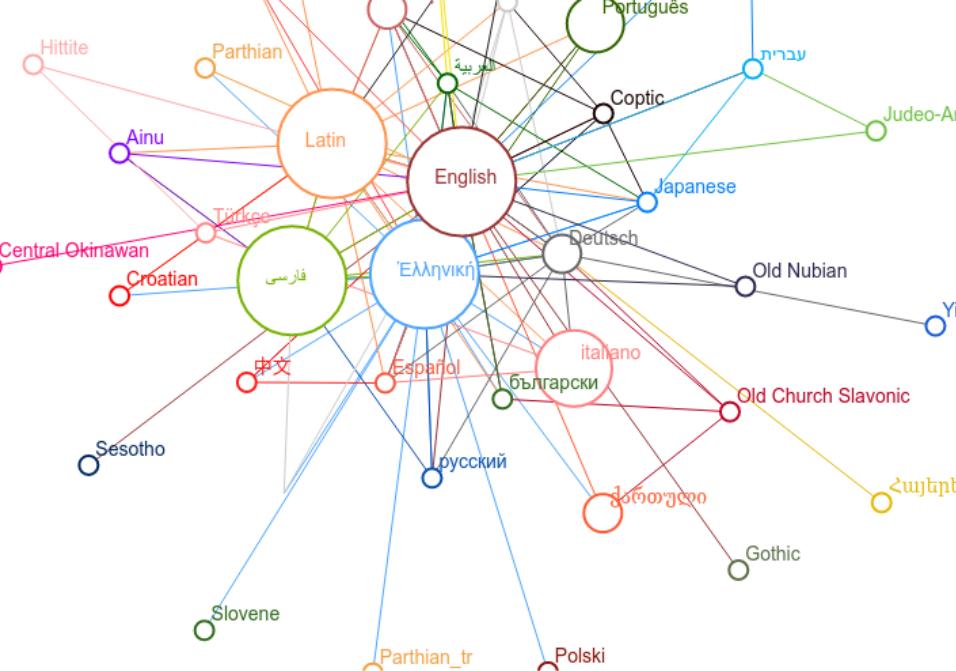
You may need some serious hardware.



Problems ?

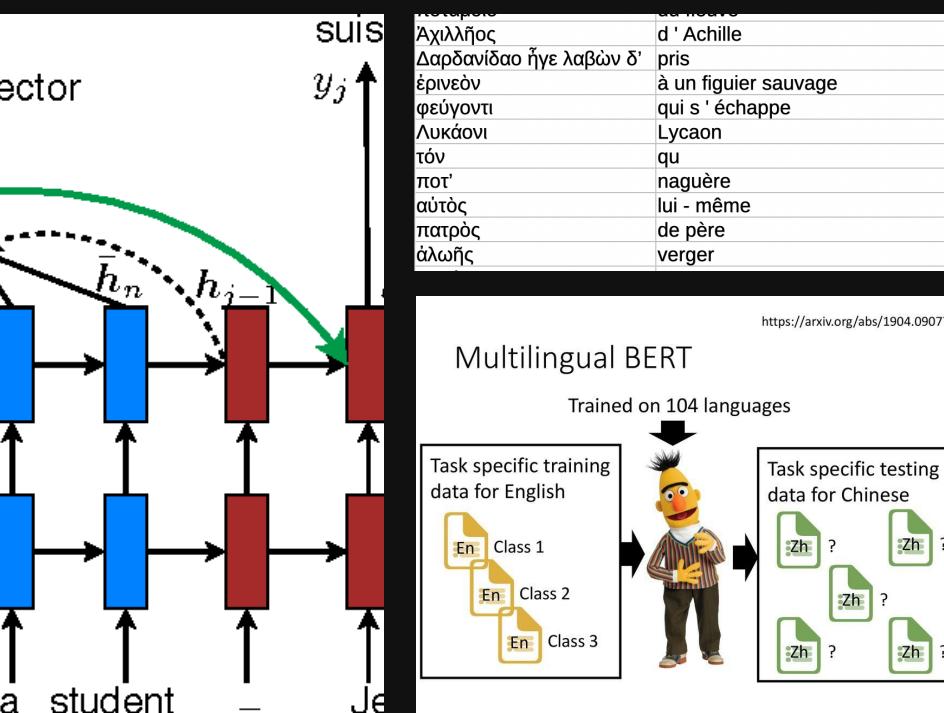
But the main obstacle is :

- very few data
- very few representative data (apart from the classical period)



MULTILINGUAL

multilingual LLMs



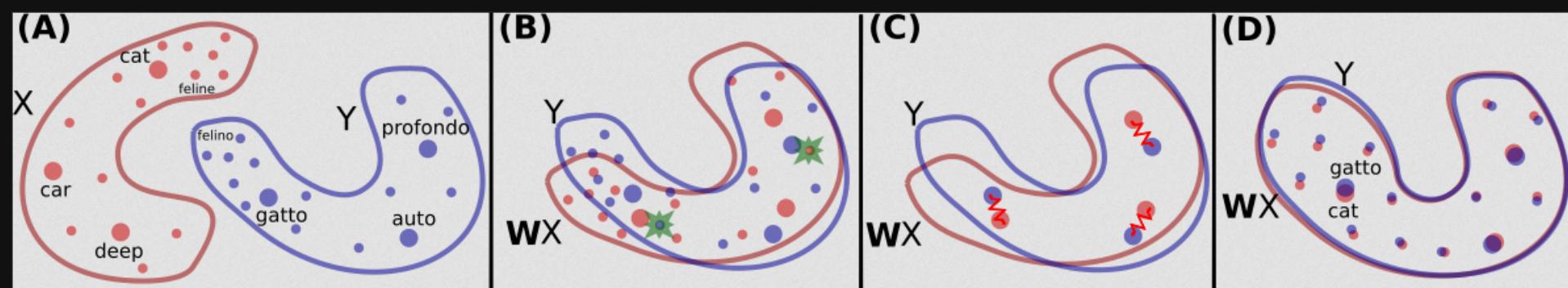
Even bigger problems

Why is it so complex ?

wmt19_translate/de-fr

- **Description** de la configuration : ensemble de données de tâche de traduction de-en WMT 2019.
- **Taille du téléchargement** : 9.71 GiB

Solution 1 : mapping monolingual spaces

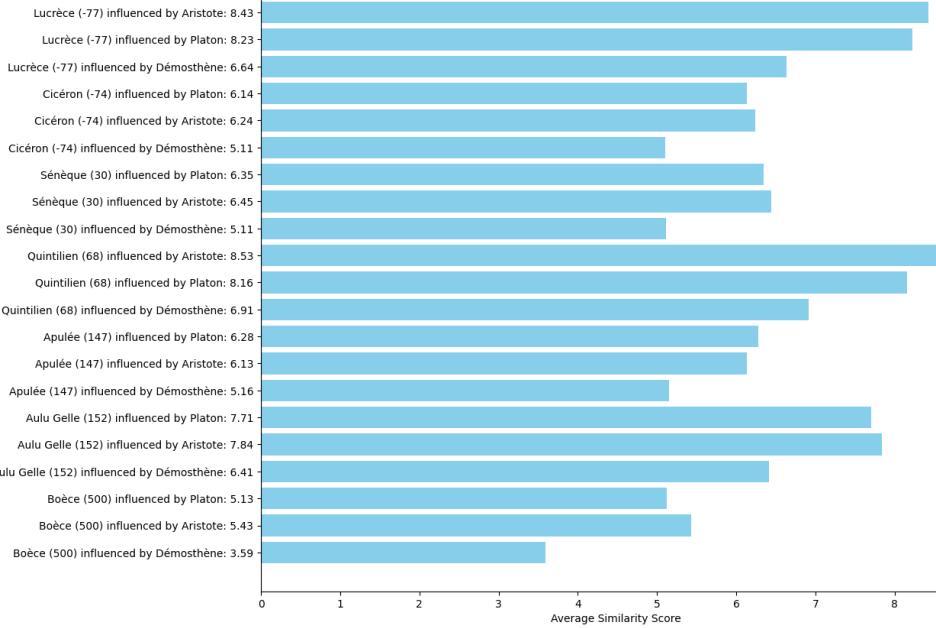


Token similarity with Multilingual BERT



```
greek_word = "λόγος"  
greek_embedding = greek_filtered_embeddings[greek_word]  
  
top_latin_words_with_scores = find_most_similar(greek_embedding, latin_filtered_embeddings)  
  
for word, score in top_latin_words_with_scores:  
    print(f"{word}: {score:.4f}")
```

```
oratio: 0.9894  
τῷ: 0.9892  
narratio: 0.9885  
καὶ: 0.9878  
liber: 0.9877  
finis: 0.9877  
scribo: 0.9877  
loquor: 0.9875  
ratio: 0.9875  
argumentum: 0.9872
```



```
_grecs = [
    "halès", "tlg1705", -624, -546),
    "ythagore", "tlg0632", -570, -495),
    "éracrite", "tlg0626", -535, -475),
    "arménide", "tlg1562", -515, -450),
    "cidamas", "tlg0610", -450, -400),
    "ntisthène", "tlg0591", -445, -365),
    "laton", "tlg0059", -427, -347),
    "émocrète", "tlg1304", -460, -370),
    "ristote", "tlg0086", -384, -322),
    "ogène de Sinope", "tlg1325", -413, -322),
    "orgias de Léontium", "tlg0593", -480, -384),
    "émosthène", "tlg0014", -384, -322),
    "pictète", "tlg0557", 50, 135)
```

```
_latins = [
    "ciceron", "phi0474", -106, -43),
    "ucrèce", "phi0550", -99, -55),
    "arron", "phi0684", -116, -27),
    "énèque", "phi1017", -4, 65),
    "énèque", "stoa0255", -4, 65),
    "énèque", "phi1014", -4, 65),
    "quintilien", "phi1002", 35, 100),
    "apulée", "phi1212", 124, 170),
    "Aulu Gelle", "phi1254", 125, 180),
    "boèce", "stoa0058", 477, 524),
    "sonius Rufus", "tlg0628", 20, 101),
    "arc Aurèle", "tlg0562", 121, 180)
```

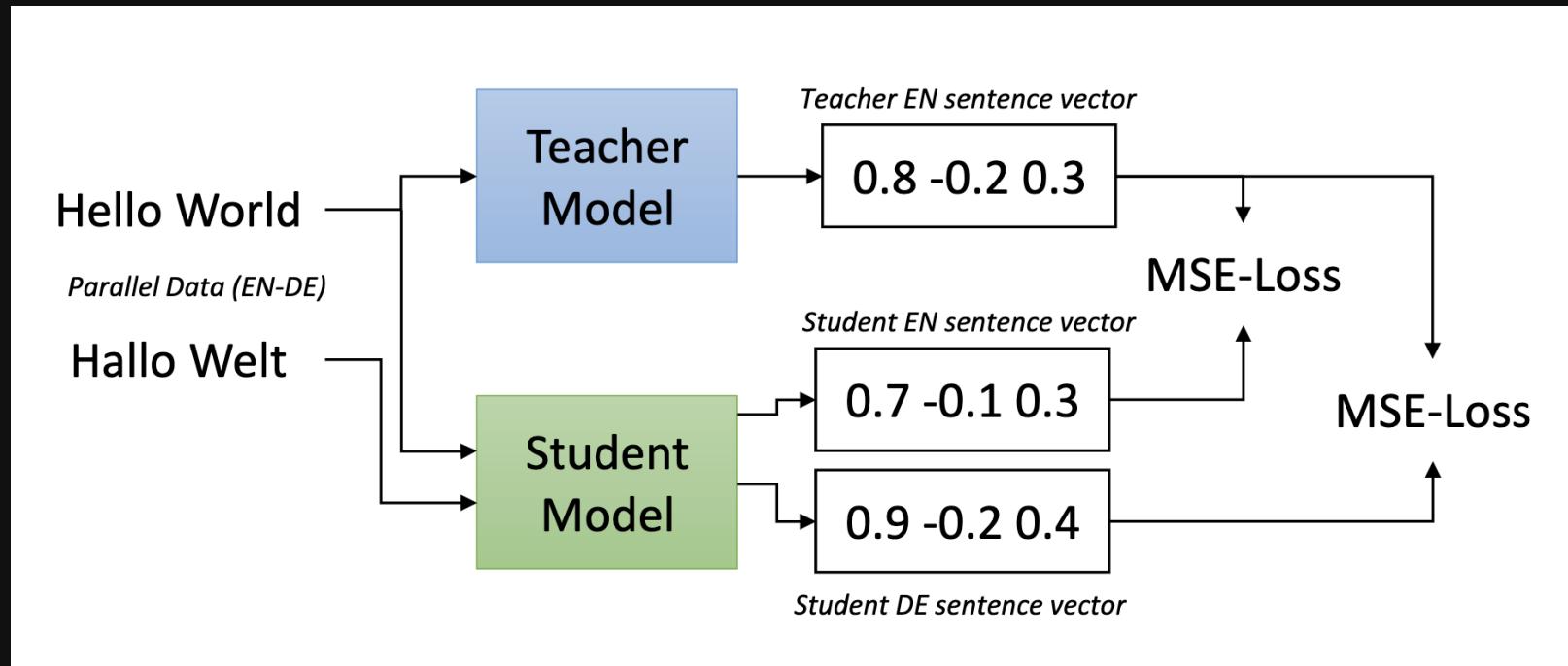
file: tlg0086.tlg025.perseus-grc2.xml, Latin File: phi0550.phi001.perseus-latin2.xml
Phrase: δῆλον δὲ καὶ ἐκ τῶν τοιούτων λόγων.
Similarity Score: 0.9494874477386475
Phrase: οὐδὲ μὲν οὖν τούτων δεδηλώται καὶ πρότερον.
Similarity Score: 0.9490494132041931
Phrase: οὐδὲ ἀμορία καὶ ἔτερα οὐδὲν.
Similarity Score: 0.948077917098999
Phrase: id licet hinc quamvis habeti cognoscere corde.
Similarity Score: 0.9477147459983826

8.12008
Mutation and Language
In vicorem cepit. Detecting Latin Allusions to Ancient Greek Texts
Anette Frank
A pivotal role in Classical Philology, with Latin authors frequently referencing Ancient Greek texts. However, automated detection of such textual references has been constrained to monolingual approaches, seeking parallels solely within the same language. SPHILBERTa is a multilingual Sentence-RoBERTa model tailored for Classical Philology, which excels at cross-lingual detection of allusions across Ancient Greek, Latin, and English. We generate new training data by automatically extracting allusions from Latin authors. We present a case study, demonstrating SPHILBERTa's capability to facilitate automated detection of Latin allusions to Greek texts. The code is available at <https://github.com/CS-CL/SPHILBERTa>.
arXiv preprint at the First Workshop on Ancient Language Processing (ALP) 2023; 9 pages, 5 tables
Language (cs.CL)
arXiv version 1
[cs.CL] for this version
50/arXiv.2308.12008

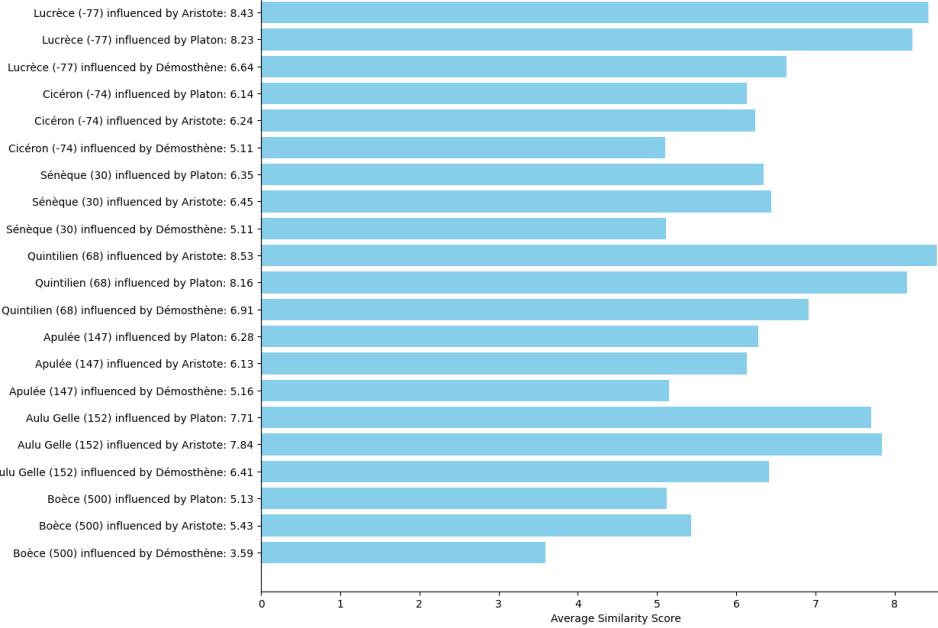
Experiment

Can we quantify proximity between philosophers ?

Model selection : sPhilberta



Nils Reimers and Iryna Gurevych. 2020. [Making Monolingual Sentence Embeddings Multilingual using Knowledge Distillation](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4512–4525, Online. Association for Computational Linguistics.



Steps

from topoi to original ideas

```
_grecs = [
    "halès", "tlg1705", -624, -546),
    "ythagore", "tlg0632", -570, -495),
    "éraclite", "tlg0626", -535, -475),
    "arménide", "tlg1562", -515, -450),
    "cidamas", "tlg0610", -450, -400),
    "ntisthène", "tlg0591", -445, -365),
    "laton", "tlg0059", -427, -347),
    "émocrите", "tlg1304", -460, -370),
    "ristote", "tlg0086", -384, -322),
    "ogène de Sinope", "tlg1325", -413, -322),
    "rgias de Léontium", "tlg0593", -480, -384),
    "émosthène", "tlg0014", -384, -322),
    "pictète", "tlg0557", 50, 135)
```

```
_latins = [
    "ciceron", "phi0474", -106, -43),
    "ucrèce", "phi0550", -99, -55),
    "arron", "phi0684", -116, -27),
    "énèque", "phi1017", -4, 65),
    "énèque", "stoa0255", -4, 65),
    "énèque", "phi1014", -4, 65),
    "quintilien", "phi1002", 35, 100),
    "apulée", "phi1212", 124, 170),
    "Aulu Gelle", "phi1254", 125, 180),
    "boèce", "stoa0058", 477, 524),
    "sonius Rufus", "tlg0628", 20, 101),
    "arc Aurèle", "tlg0562", 121, 180)
```

file: tlg0086.tlg025.perseus-grc2.xml, Latin File: phi0550.phi001.perseus-latin
Phrase: δῆλον δὲ καὶ ἐκ τῶν τοιούτων λόγων.
Phrase: id quod iam supra tibi paulo ostendimus ante.
Similarity Score: 0.9494874477386475
Phrase: μερὶ μὲν οὖν τούτων δεδηλώται καὶ πρότερον.
Phrase: quo magis incemptum pergam pertexere dictis.
Similarity Score: 0.9490494132041931
Phrase: ξοτὶ δὲ ἀμορία καὶ ἔτερα μερὶ αὐτῶν.
Phrase: quae tibi posterius largo sermone probabo.
Similarity Score: 0.9480779170989999
Phrase: μερὶ δὲ οὐδὲν ήττον συμβαίνει τὰ αὐτὰ ἀμορεῖν.
Phrase: id licet hinc quamvis habeti cognoscere corde.
Similarity Score: 0.9477147459983826

8.12008
Mutation and Language
in victorem cepit. Detecting Latin Allusions to Ancient Greek Topoi
Anette Frank
pivot role in Classical Philology, with Latin authors frequently referencing Ancient Greek texts. However, this cross-lingual reference detection has been constrained to monolingual approaches, seeking parallels solely within one language. We introduce SPhILBERTa, a multilingual Sentence-RoBERTa model tailored for Classical Philology, which excels at cross-lingual reference detection across Ancient Greek, Latin, and English. We generate new training data by automatically extracting allusions from Latin authors. We present a case study, demonstrating SPhILBERTa's capability to facilitate automated detection of Latin allusions to Greek topoi. This work is available at <https://arxiv.org/abs/2308.12008>.
This work was presented at the First Workshop on Ancient Language Processing (ALP) 2023; 9 pages, 5 tables
Language (cs.CL)
[cs.CL] [cs.LG]
[cs.CL] for this version
arXiv:2308.12008

Step 1 : corpus selection

Global open
source data
(First1KGreek,
Perseus)

Thales
Pythagoras
Heraclitus
Parmenides
Alcidamas
Antisthenes
Plato
Democritus
Aristotles
Diogenes of sinope
Gorgias
Demosthenes
Epictetus
~ 50 000 unique
sentences



Cicero
Lucretius
Varro
Seneca
Quintilian
Apuleius
Aulus Gellius
Boetius
Musonius Rufus
Marcus Aurelius
~ 100 000 unique
sentences

Step 2 : topoi versus peculiarities

Thales
 Pythagoras
 Heraclitus
 Parmenides
 Alcidamas
 Antisthenes
 Plato
 Democritus
 Aristotles
 Diogenes of sinope
 Gorgias
 Demosthenes
 Epictetus
 ~ 50 000 unique sentences

Docs	Word 1	Word 2	Word 3	Author
Doc1	3	0	1	A
Doc2	0	4	1	B

$$E_{ij} = (N_i \times N_j) / N$$

N_i : all words in class i .

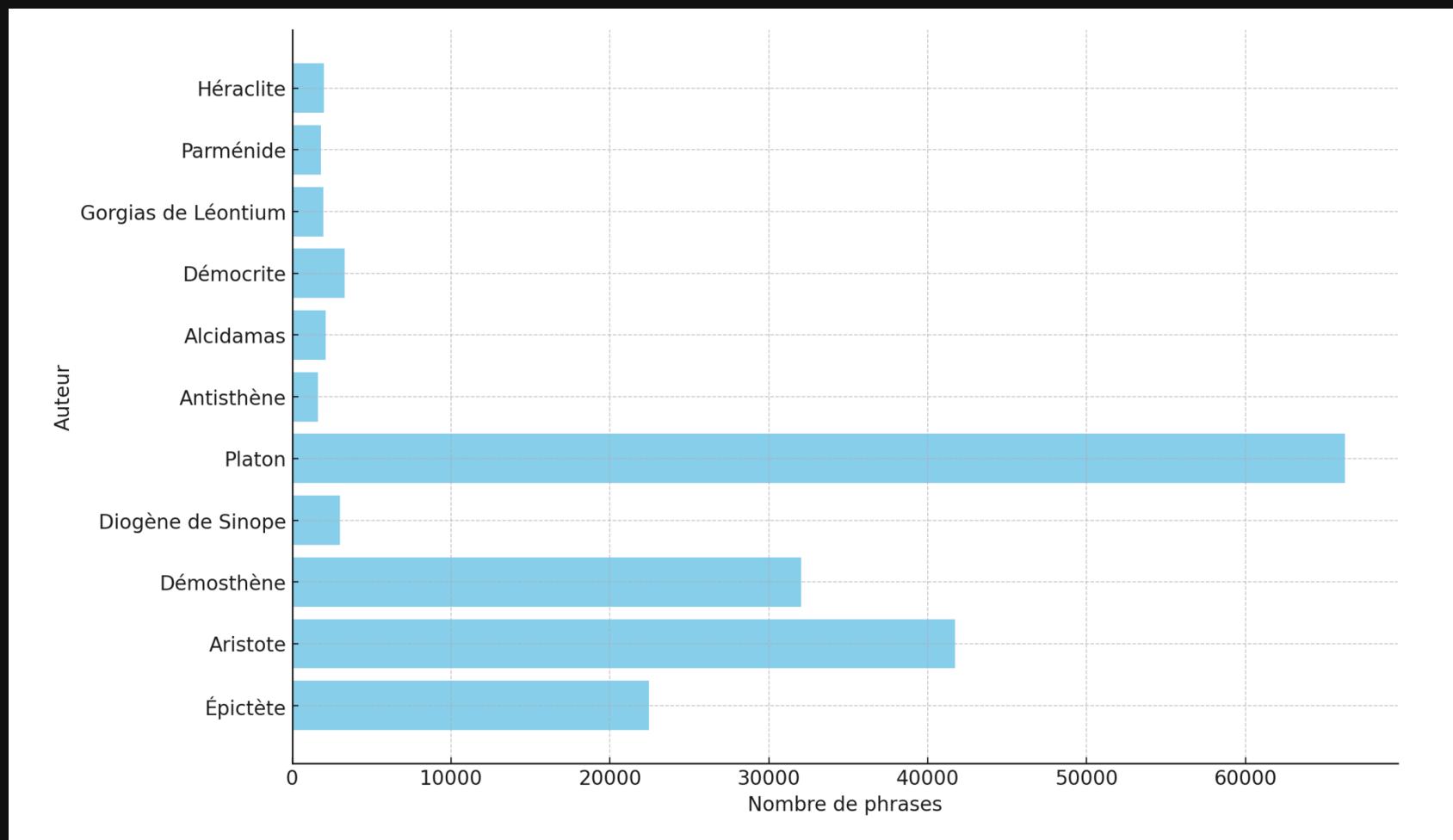
N_j : all words j in all classes.

N : all words in all classes.

- Word 1 pour Auteur A : $E = (4 * 3) / 9 = 1.33$.
- Word 1 pour Auteur B : $E = (5 * 3) / 9 = 1.67$.

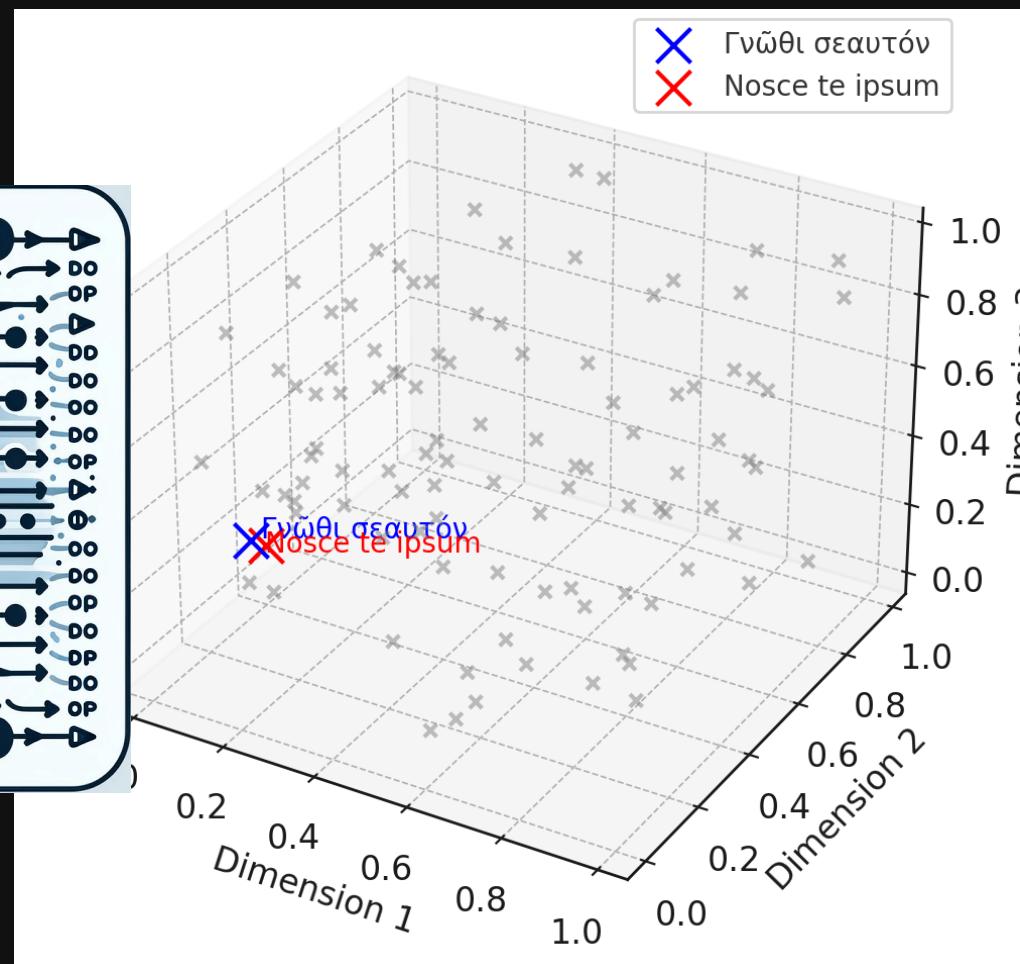
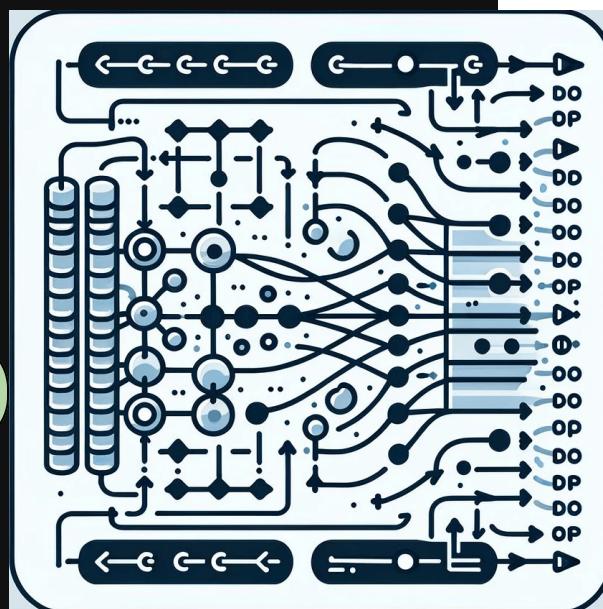
$$\begin{aligned} \bullet \quad \chi^2 &= \frac{(3-1.33)^2}{1.33} + \frac{(0-1.67)^2}{1.67} = \frac{2.77}{1.33} + \frac{2.79}{1.67} = 2.08 + 1.67 = 3.75 \\ \bullet \quad \chi^2 &= \frac{(0-1.78)^2}{1.78} + \frac{(4-2.22)^2}{2.22} = \frac{3.17}{1.78} + \frac{3.17}{2.22} = 1.78 + 1.43 = 3.21 \\ \bullet \quad \chi^2 &= \frac{(1-0.89)^2}{0.89} + \frac{(1-1.11)^2}{1.11} = \frac{0.0121}{0.89} + \frac{0.0121}{1.11} = 0.014 + 0.011 = 0.025 \end{aligned}$$

Step 3 : filtering corpus



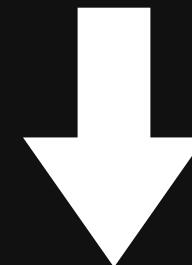
Step 4 : encoding

latin sentence
greek sentence



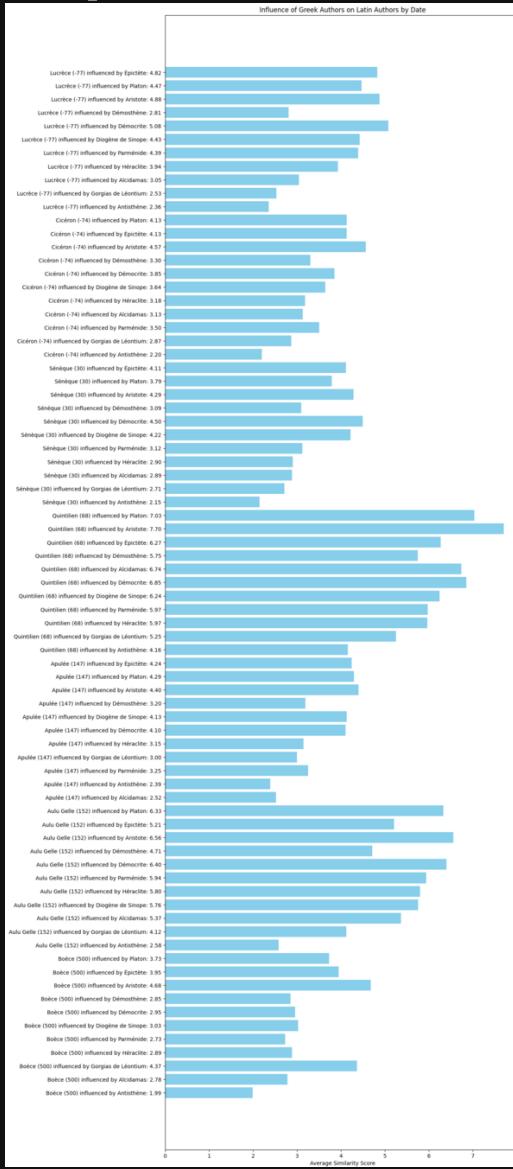
Step 4 : raw results and re-calculated

```
Distribution des similarités filtrées par auteur grec :  
Auteur : tlg0014, Nombre de similarités : 1251200  
Auteur : tlg0086, Nombre de similarités : 1669362  
Auteur : tlg0557, Nombre de similarités : 824366  
Auteur : tlg0059, Nombre de similarités : 3359990  
Auteur : tlg1562, Nombre de similarités : 15828  
Auteur : tlg0626, Nombre de similarités : 14775  
Auteur : tlg1304, Nombre de similarités : 56673  
Auteur : tlg0610, Nombre de similarités : 23685  
Auteur : tlg1325, Nombre de similarités : 38615  
Auteur : tlg0593, Nombre de similarités : 13377  
Auteur : tlg0591, Nombre de similarités : 3633
```



Similarity means / number of sentences per
author (total)

Step 5 : final results



Main results :

Aristotle, although with less sentences, is much more represented in latin authors on the whole period

Plato remains one of the main inspirations through time

Some of the smallest authors (in terms of frequency) are still represented, mostly in earlier latin authors

Epictetes is definitely influenced by contemporary latin authors

Quintilian seems to be the author where global similarities are the most referenced

Contact

Thank you for your attention

Github link for experiment

marianne.reboul@ens-lyon.fr

github : OdysseusPolymetis

Useful references

- Yousef, T., Palladino, C., Shamsian, F., d'Orange Ferreira, A., & Ferreira dos Reis, M. (2022). An automatic model and Gold Standard for translation alignment of Ancient Greek. In Proceedings of the Language Resources and Evaluation Conference (pp. 5894–5905). European Language Resources Association. Marseille, France. Retrieved from <https://aclanthology.org/2022.lrec-1.634>
- Nils Reimers and Iryna Gurevych. 2020. Making Monolingual Sentence Embeddings Multilingual using Knowledge Distillation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4512–4525, Online. Association for Computational Linguistics.
- Riemenschneider, F., & Frank, A. (2023). Graecia capta ferum victorem cepit. Detecting Latin Allusions to Ancient Greek Literature. In *Proceedings of the First Workshop on Ancient Language Processing*. Association for Computational Linguistics. Varna, Bulgaria. Retrieved from <https://arxiv.org/abs/2308.12008>
- Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- MZES Data Lab. (n.d.). BERT and explainable AI. Retrieved from <https://www.mzes.uni-mannheim.de/socialsciencedatalab/article/bert-explainable-ai/>

Contact

Thank you for your attention

Colab link for experiment

marianne.reboul@ens-lyon.fr

[github : OdysseusPolymetis](#)

Bibliographie

- Yousef, T., Palladino, C., & Jänicke, S. (2022). Transformer-based named entity recognition for ancient greek.
- Yousef, T., Palladino, C., Shamsian, F., Ferreira, A. D. O., & dos Reis, M. F. (2022, June). An automatic model and Gold Standard for translation alignment of Ancient Greek. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference* (pp. 5894-5905).
- Riemenschneider, F., & Frank, A. (2023). Exploring large language models for classical philology. *arXiv preprint arXiv:2305.13698*.
- Artetxe, M., Labaka, G., & Agirre, E. (2018). A robust self-learning method for fully unsupervised cross-lingual mappings of word embeddings. *arXiv preprint arXiv:1805.06297*.
- Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.