# LANGSAMP: Language-Script Aware Multilingual Pretraining

Yihong Liu<sup>1,2,\*</sup>, Haotian Ye<sup>1,2,\*</sup>, Chunlan Ma<sup>1,2</sup>, Mingyang Wang<sup>1,2,3</sup>, Hinrich Schütze<sup>1,2</sup>







Munich Center for Machine Learning

<sup>1</sup>Center for Information and Language Processing, LMU Munich <sup>2</sup>Munich Center for Machine Learning (MCML) <sup>3</sup>Bosch Center for Artificial Intelligence

September 24, 2025

#### Outline

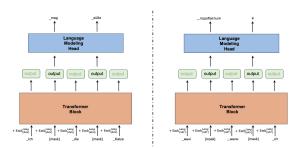


1 LangSAMP

2 Experiments

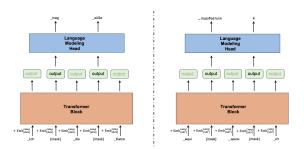
3 Analysis





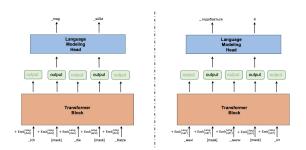
• Early crosslingual models like XLM leverage language embeddings.





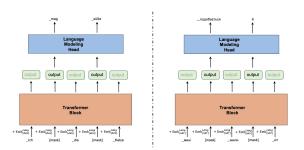
- Early crosslingual models like XLM leverage language embeddings.
  - learnable vectors





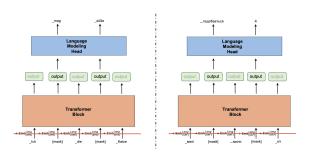
- Early crosslingual models like XLM leverage language embeddings.
  - learnable vectors
  - capture language-specific information





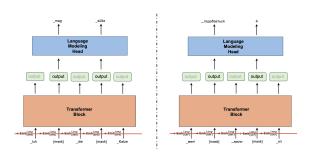
- Early crosslingual models like XLM leverage language embeddings.
  - learnable vectors
  - capture language-specific information
  - useful for guiding generation, e.g., MT





Recent mPLMs like XLM-R remove such language embeddings.

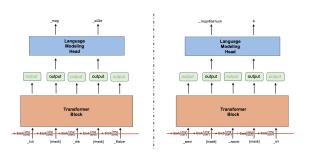




Recent mPLMs like XLM-R remove such language embeddings.

• What's bad about such removal?

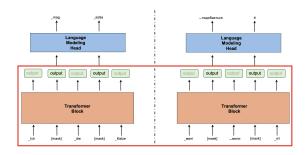




Recent mPLMs like XLM-R remove such language embeddings.

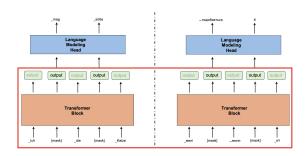
- What's bad about such removal?
  - The contextual token embeddings have to encode all language-specific information
  - This may hinder language neutrality
  - Language neutrality is usually important for multilinguality, e.g., (zero-shot) crosslingual transfer.





• What's good about such removal?





- What's good about such removal?
  - Universal text encoder without requiring language IDs as input
  - The backbone model can be used effectively for downstream tasks



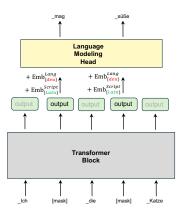
- Why not integrate language/script embeddings wisely so that
  - representations' language-neutrality is improved
  - the backbone remains the same as common mPLMs

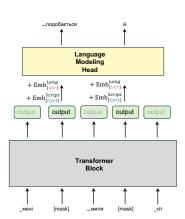
7

### LANGSAMP- Modeling

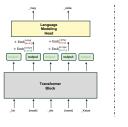


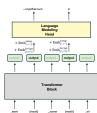
We add the language and script embedding to the outputs of the Transformer blocks at token position i:  $o_i = h_i + E_s^{Lang} + E_s^{Script}$ .





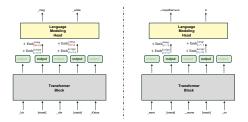






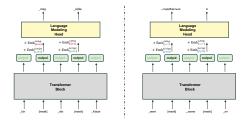


**Note**: language/script embeddings are **not required** to obtain the final Transformer output, i.e., the final contextual token embeddings.



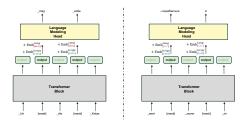
• No language or script IDs as input needed in fine-tuning.





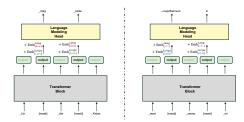
- No language or script IDs as input needed in fine-tuning.
- The backbone remains the same as most mPLMs.





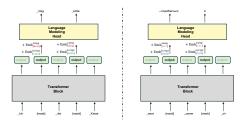
- No language or script IDs as input needed in fine-tuning.
- The backbone remains the same as most mPLMs.
- It can be fine-tuned in the standard way in NLP pipelines.





- No language or script IDs as input needed in fine-tuning.
- The backbone remains the same as most mPLMs.
- It can be fine-tuned in the standard way in NLP pipelines.
  - Token (and position) embeddings + Transformer blocks as **backbone**.





- No language or script IDs as input needed in fine-tuning.
- The backbone remains the same as most mPLMs.
- It can be fine-tuned in the standard way in NLP pipelines.
  - Token (and position) embeddings + Transformer blocks as **backbone**.
  - A task-specific classifier is added on top of it.

#### Outline



1 LANGSAMP

Experiments

3 Analysis



- Continued pretraining XLM-R on Glot500-c
- head: the language covered by XLM-R.
- tail: the language not covered by XLM-R

		tail		head	1	Latn	no	n-Latn		all
	Baseline	LANGSAMP								
SR-B	36.9 (0.0)	<b>39.5</b> (0.0)	60.6 (0.0)	61.3 (0.0)	40.7 (0.0)	42.8 (0.0)	51.2 (0.0)	<b>53.5</b> (0.0)	42.9 (0.0)	<b>45.1</b> (0.0)
SR-T	56.9 (0.0)	<b>58.6</b> (0.0)	74.8 (0.0)	<b>76.1</b> (0.0)	67.5 (0.0)	<b>68.7</b> (0.0)	73.7 (0.0)	<b>75.6</b> (0.0)	69.7 (0.0)	<b>71.1</b> (0.0)
Taxi1500	47.1 (4.8)	50.8 (2.4)	59.9 (2.9)	<b>61.2</b> (1.2)	48.2 (4.6)	<b>51.7</b> (2.1)	58.8 (3.1)	60.1 (1.7)	50.3 (4.2)	<b>53.4</b> (2.0)
SIB200	69.0 (1.4)	70.2 (1.9)	82.2 (1.4)	82.6 (1.2)	72.1 (1.3)	73.1 (1.8)	81.1 (1.5)	81.7 (1.2)	75.0 (1.3)	<b>75.9</b> (1.6)
NER	60.1 (0.6)	<b>60.8</b> (0.8)	64.0 (0.6)	<b>64.1</b> (0.6)	67.0 (0.5)	<b>67.6</b> (0.6)	53.9 (0.7)	<b>53.9</b> (0.5)	62.2 (0.5)	<b>62.6</b> (0.6)
POS	61.3 (1.0)	61.4 (0.9)	76.0 (0.4)	<b>76.2</b> (0.4)	74.6 (0.5)	74.5 (0.4)	66.2 (1.0)	<b>66.8</b> (0.8)	71.5 (0.6)	<b>71.6</b> (0.5)



- Continued pretraining XLM-R on Glot500-c
- head: the language covered by XLM-R.
- tail: the language not covered by XLM-R

		tail		head	I	Latn	no	n-Latn		all
	Baseline	LANGSAMP	Baseline	LANGSAMP	Baseline	LANGSAMP	Baseline	LANGSAMP	Baseline	LANGSAMP
SR-B	36.9 (0.0)	<b>39.5</b> (0.0)	60.6 (0.0)	<b>61.3</b> (0.0)	40.7 (0.0)	<b>42.8</b> (0.0)	51.2 (0.0)	<b>53.5</b> (0.0)	42.9 (0.0)	<b>45.1</b> (0.0)
SR-T	56.9 (0.0)	<b>58.6</b> (0.0)	74.8 (0.0)	<b>76.1</b> (0.0)	67.5 (0.0)	68.7 (0.0)	73.7 (0.0)	<b>75.6</b> (0.0)	69.7 (0.0)	<b>71.1</b> (0.0)
Taxi1500	47.1 (4.8)	50.8 (2.4)	59.9 (2.9)	<b>61.2</b> (1.2)	48.2 (4.6)	<b>51.7</b> (2.1)	58.8 (3.1)	60.1 (1.7)	50.3 (4.2)	<b>53.4</b> (2.0)
SIB200	69.0 (1.4)	70.2 (1.9)	82.2 (1.4)	82.6 (1.2)	72.1 (1.3)	<b>73.1</b> (1.8)	81.1 (1.5)	81.7 (1.2)	75.0 (1.3)	<b>75.9</b> (1.6)
NER	60.1 (0.6)	<b>60.8</b> (0.8)	64.0 (0.6)	<b>64.1</b> (0.6)	67.0 (0.5)	<b>67.6</b> (0.6)	<b>53.9</b> (0.7)	<b>53.9</b> (0.5)	62.2 (0.5)	<b>62.6</b> (0.6)
POS	61.3 (1.0)	<b>61.4</b> (0.9)	76.0 (0.4)	<b>76.2</b> (0.4)	<b>74.6</b> (0.5)	74.5 (0.4)	66.2 (1.0)	<b>66.8</b> (0.8)	71.5 (0.6)	<b>71.6</b> (0.5)

Consistently better performance on both head and tail languages.



- Continued pretraining XLM-R on Glot500-c
- head: the language covered by XLM-R.
- tail: the language not covered by XLM-R

		tail		head		Latn	no	n-Latn	all		
	Baseline	LANGSAMP	Baseline	LANGSAMP	Baseline	LANGSAMP	Baseline	LANGSAMP	Baseline	LANGSAMP	
SR-B	36.9 (0.0)	<b>39.5</b> (0.0)	60.6 (0.0)	61.3 (0.0)	40.7 (0.0)	42.8 (0.0)	51.2 (0.0)	<b>53.5</b> (0.0)	42.9 (0.0)	<b>45.1</b> (0.0)	
SR-T	56.9 (0.0)	<b>58.6</b> (0.0)	74.8 (0.0)	<b>76.1</b> (0.0)	67.5 (0.0)	<b>68.7</b> (0.0)	73.7 (0.0)	<b>75.6</b> (0.0)	69.7 (0.0)	71.1 (0.0)	
Taxi1500	47.1 (4.8)	50.8 (2.4)	59.9 (2.9)	<b>61.2</b> (1.2)	48.2 (4.6)	51.7 (2.1)	58.8 (3.1)	60.1 (1.7)	50.3 (4.2)	<b>53.4</b> (2.0)	
SIB200	69.0 (1.4)	70.2 (1.9)	82.2 (1.4)	82.6 (1.2)	72.1 (1.3)	73.1 (1.8)	81.1 (1.5)	81.7 (1.2)	75.0 (1.3)	<b>75.9</b> (1.6)	
NER	60.1 (0.6)	<b>60.8</b> (0.8)	64.0 (0.6)	<b>64.1</b> (0.6)	67.0 (0.5)	<b>67.6</b> (0.6)	<b>53.9</b> (0.7)	<b>53.9</b> (0.5)	62.2 (0.5)	<b>62.6</b> (0.6)	
POS	61.3 (1.0)	<b>61.4</b> (0.9)	76.0 (0.4)	<b>76.2</b> (0.4)	<b>74.6</b> (0.5)	74.5 (0.4)	66.2 (1.0)	<b>66.8</b> (0.8)	71.5 (0.6)	<b>71.6</b> (0.5)	

- Consistently better performance on both head and tail languages.
- Both non-Latin and Latin languages benefit from LANGSAMP.



- Continued pretraining XLM-R on Glot500-c
- head: the language covered by XLM-R.
- tail: the language not covered by XLM-R

		tail		head		Latn	no	n-Latn	all		
	Baseline	LANGSAMP	Baseline	LANGSAMP	Baseline	LANGSAMP	Baseline	LANGSAMP	Baseline	LANGSAMP	
SR-B	36.9 (0.0)	<b>39.5</b> (0.0)	60.6 (0.0)	61.3 (0.0)	40.7 (0.0)	42.8 (0.0)	51.2 (0.0)	<b>53.5</b> (0.0)	42.9 (0.0)	<b>45.1</b> (0.0)	
SR-T	56.9 (0.0)	<b>58.6</b> (0.0)	74.8 (0.0)	<b>76.1</b> (0.0)	67.5 (0.0)	<b>68.7</b> (0.0)	73.7 (0.0)	<b>75.6</b> (0.0)	69.7 (0.0)	71.1 (0.0)	
Taxi1500	47.1 (4.8)	50.8 (2.4)	59.9 (2.9)	<b>61.2</b> (1.2)	48.2 (4.6)	51.7 (2.1)	58.8 (3.1)	60.1 (1.7)	50.3 (4.2)	<b>53.4</b> (2.0)	
SIB200	69.0 (1.4)	70.2 (1.9)	82.2 (1.4)	82.6 (1.2)	72.1 (1.3)	73.1 (1.8)	81.1 (1.5)	81.7 (1.2)	75.0 (1.3)	<b>75.9</b> (1.6)	
NER	60.1 (0.6)	<b>60.8</b> (0.8)	64.0 (0.6)	<b>64.1</b> (0.6)	67.0 (0.5)	<b>67.6</b> (0.6)	<b>53.9</b> (0.7)	<b>53.9</b> (0.5)	62.2 (0.5)	<b>62.6</b> (0.6)	
POS	61.3 (1.0)	<b>61.4</b> (0.9)	76.0 (0.4)	<b>76.2</b> (0.4)	<b>74.6</b> (0.5)	74.5 (0.4)	66.2 (1.0)	<b>66.8</b> (0.8)	71.5 (0.6)	<b>71.6</b> (0.5)	

- Consistently better performance on both head and tail languages.
- Both non-Latin and Latin languages benefit from LANGSAMP.
- Improvements can vary slightly across different task types.

#### Outline



1 LANGSAMP

2 Experiments

Analysis

### Ablation Study



		SR-B			SR-T		Т	axi150	00	Ş	SIB200		NER			POS		
	tail	head	all	tail	head	all	tail	head	all	tail	head	all	tail	head	all	tail	head	all
vanilla model	11.9	56.4	23.2	46.0	77.7	68.6	18.1	58.6	28.4	56.1	83.0	68.3	55.1	62.8	59.3	49.9	75.7	67.8
w/ <b>E</b> <sup>Lang</sup> w/ <b>E</b> <sup>Script</sup>	13.1	<u>57.9</u>	24.5	49.1	79.0	70.5	18.3	58.5	28.5	57.2	82.7	68.8	55.2	63.0	59.5	49.9	<u>75.8</u>	67.8
w/ <b>E</b> <sup>Script</sup>	12.5	57.4	23.9	48.3	78.4	69.8	18.5	57.0	28.2	56.6	82.1	68.2	55.1	62.4	59.0	50.8	76.2	68.4
w/ $\boldsymbol{E}^{Lang}$ and $\boldsymbol{E}^{Script}$	13.4	58.7	24.9	49.1	79.5	70.8	20.6	58.8	30.3	57.9	83.0	69.3	54.9	61.6	58.6	49.7	75.6	67.6

### Ablation Study



	SR-B			SR-T			Taxi1500			SIB200			NER			POS		
	tail	head	all	tail	head	all	tail	head	all	tail	head	all	tail	head	all	tail	head	all
vanilla model	11.9	56.4	23.2	46.0	77.7	68.6	18.1	58.6	28.4	56.1	83.0	68.3	55.1	62.8	59.3	49.9	75.7	67.8
w/ <b>E</b> <sup>Lang</sup> w/ <b>E</b> <sup>Script</sup>	13.1	<u>57.9</u>	24.5	49.1	79.0	70.5	18.3	58.5	28.5	57.2	82.7	68.8	55.2	63.0	59.5	49.9	<u>75.8</u>	67.8
											82.1							
w/ $\boldsymbol{E}^{Lang}$ and $\boldsymbol{E}^{Script}$	13.4	58.7	24.9	49.1	79.5	70.8	20.6	58.8	30.3	57.9	83.0	69.3	54.9	61.6	58.6	49.7	75.6	67.6

• Both language and script embeddings help.

### Ablation Study

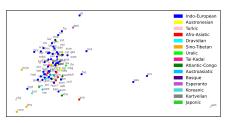


	SR-B			SR-T			Taxi1500			SIB200			NER			POS		
	tail	head	all	tail	head	all	tail	head	all	tail	head	all	tail	head	all	tail	head	all
vanilla model	11.9	56.4	23.2	46.0	77.7	68.6	18.1	58.6	28.4	56.1	83.0	68.3	55.1	62.8	59.3	49.9	75.7	67.8
w/ <b>E</b> <sup>Lang</sup> w/ <b>E</b> <sup>Script</sup>	13.1	<u>57.9</u>	24.5	49.1	79.0	70.5	18.3	58.5	28.5	57.2	82.7	68.8	55.2	63.0	59.5	49.9	<u>75.8</u>	67.8
											82.1							
w/ $\boldsymbol{E}^{Lang}$ and $\boldsymbol{E}^{Script}$	13.4	58.7	24.9	49.1	79.5	70.8	20.6	58.8	30.3	57.9	83.0	69.3	54.9	61.6	58.6	49.7	75.6	67.6

- Both language and script embeddings help.
- Improvement varies across task types.

### Visualization of Auxiliary Embeddings with PCA

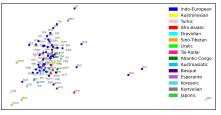






### Visualization of Auxiliary Embeddings with PCA



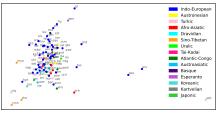




• Similar or related languages/scripts are located close to each other.

### Visualization of Auxiliary Embeddings with PCA





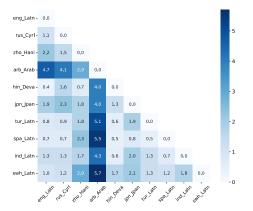


- Similar or related languages/scripts are located close to each other.
- This shows that the auxiliary embeddings capture language- and script-specific information.

### Similarity of Parallel Sentences Across Languages



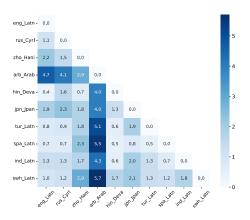
We selected 10 languages: eng\_Latn, rus\_Cyrl, zho\_Hani, arb\_Arab, hin\_Deva, jpn\_Jpan, tur\_Latn, spa\_Latn, ind\_Latn, and swa\_Latn. Pairwise cosine similarity of sentence representations is computed (baseline and LANGSAMP). Improvement (by percentage) is shown.



### Similarity of Parallel Sentences Across Languages



We selected 10 languages: eng\_Latn, rus\_Cyrl, zho\_Hani, arb\_Arab, hin\_Deva, jpn\_Jpan, tur\_Latn, spa\_Latn, ind\_Latn, and swa\_Latn. Pairwise cosine similarity of sentence representations is computed (baseline and LANGSAMP). Improvement (by percentage) is shown.

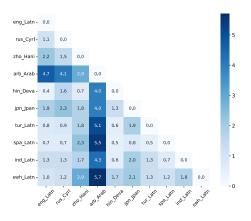


 Similarity between any two languages is improved in LANGSAMP.

### Similarity of Parallel Sentences Across Languages



We selected 10 languages: eng\_Latn, rus\_Cyrl, zho\_Hani, arb\_Arab, hin\_Deva, jpn\_Jpan, tur\_Latn, spa\_Latn, ind\_Latn, and swa\_Latn. Pairwise cosine similarity of sentence representations is computed (baseline and LANGSAMP). Improvement (by percentage) is shown.



- Similarity between any two languages is improved in LANGSAMP.
- The enhancement is especially noticeable for typologically distinct languages.

### Case Study: Source Language Selection



Instead of always using English, we select the **language** as source language whose embedding is most cosine-similar to the target language for *zero-shot crosslingual transfer*.

	tail		he	ad	La	itn	non-	Latn	all		
	English	Source									
Taxi1500	47.3	48.3	59.1	60.3	48.4	49.0	58.1	60.5	50.2	51.2	
SIB200	67.9	67.9	81.2	81.6	71.0	71.1	80.3	80.6	74.0	74.2	
NER	61.2	61.7	64.1	65.6	67.5	66.9	54.6	58.5	62.8	63.8	
POS	63.2	53.8	77.0	72.3	75.5	68.4	68.1	63.6	72.8	66.6	

### Case Study: Source Language Selection



Instead of always using English, we select the **language** as source language whose embedding is most cosine-similar to the target language for *zero-shot crosslingual transfer*.

	ta	iil	he	ad	La	tn	non-	Latn	all		
	English	Source									
Taxi1500	47.3	48.3	59.1	60.3	48.4	49.0	58.1	60.5	50.2	51.2	
SIB200	67.9	67.9	81.2	81.6	71.0	71.1	80.3	80.6	74.0	74.2	
NER	61.2	61.7	64.1	65.6	67.5	66.9	54.6	58.5	62.8	63.8	
POS	63.2	53.8	77.0	72.3	75.5	68.4	68.1	63.6	72.8	66.6	

• This shows that language embeddings facilitate the selection of optimal source languages for more effective crosslingual transfer.

#### Conclusion and Future Work



**LangSAMP:** a multilingual pretraining approach that adds auxiliary language and script embeddings *after* transformer layers to encourage language-neutral representations.

- Simple yet effective for improving language-neutrality.
- Consistent gains over baselines across downstream tasks.
- Auxiliary embeddings may be leveraged as useful byproducts.

#### **Possible Future Directions:**

- Further steer *middle-layer* representations toward neutrality using parallel data with contrastive objectives.
- Leverage embeddings for controlled multilingual generation in decoder-only or encoder-decoder models.

## Thank you for your attention!