

# Attention-Informed Mixed-Language Training for Zero-shot Cross-lingual Task-oriented Dialogue Systems

Zihan Liu, Genta Indra Winata, Zhaojiang Lin,  
Peng Xu, Pascale Fung



# Background

- Supervised neural-based approaches have shown the effectiveness for natural language processing.

# Background

- Supervised neural-based approaches have shown the effectiveness for natural language processing.
- However, they heavily rely on large amounts of training data, which makes them not scalable to low-resource languages.

# Background

- Supervised neural-based approaches have shown the effectiveness for natural language processing.
- However, they heavily rely on large amounts of training data, which makes them not scalable to low-resource languages.
- A straightforward idea is to adapt the model from the high-resource language into the low-resource languages.

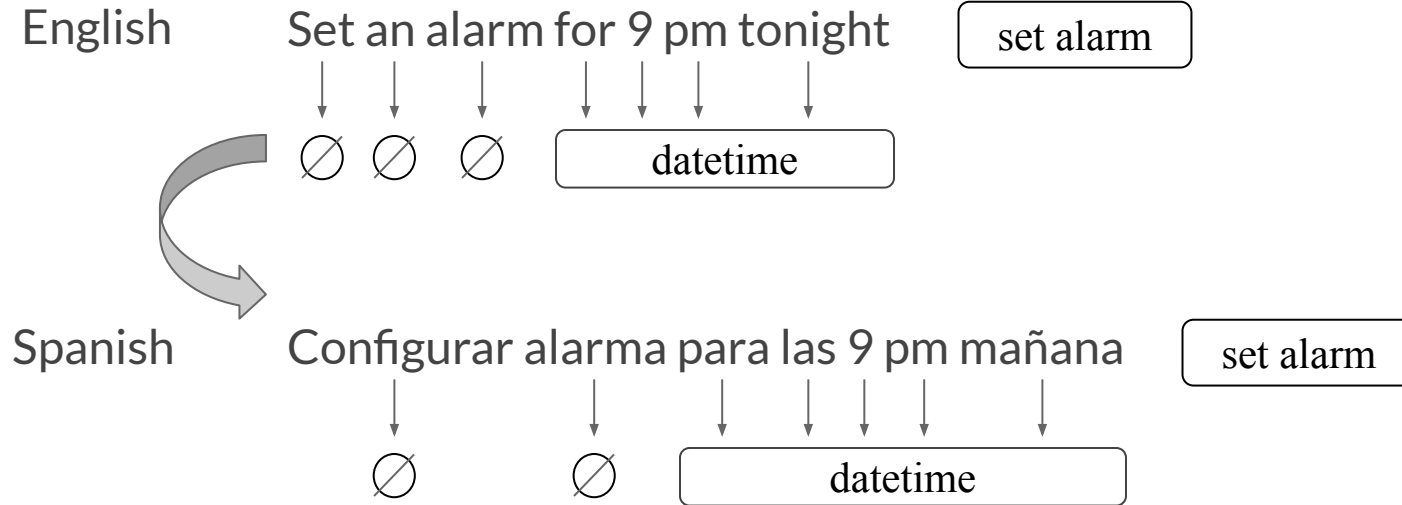
# Cross-lingual Task-oriented Dialogue Systems

- Dialogue State Tracking (DST)



# Cross-lingual Task-oriented Dialogue Systems

- Natural Language Understanding (NLU)



# Straightforward solutions

1. Translate training set from source language to target language
2. Translate test samples

# Straightforward solutions

1. Translate training set from source language to target language
2. Translate test samples

## Problems

1. We need large amounts of resources to build machine translation systems.
2. Machine translation systems perform badly if the source language and target languages are unrelated languages (e.g., English and Chinese).



# Cross-lingual Adaptation

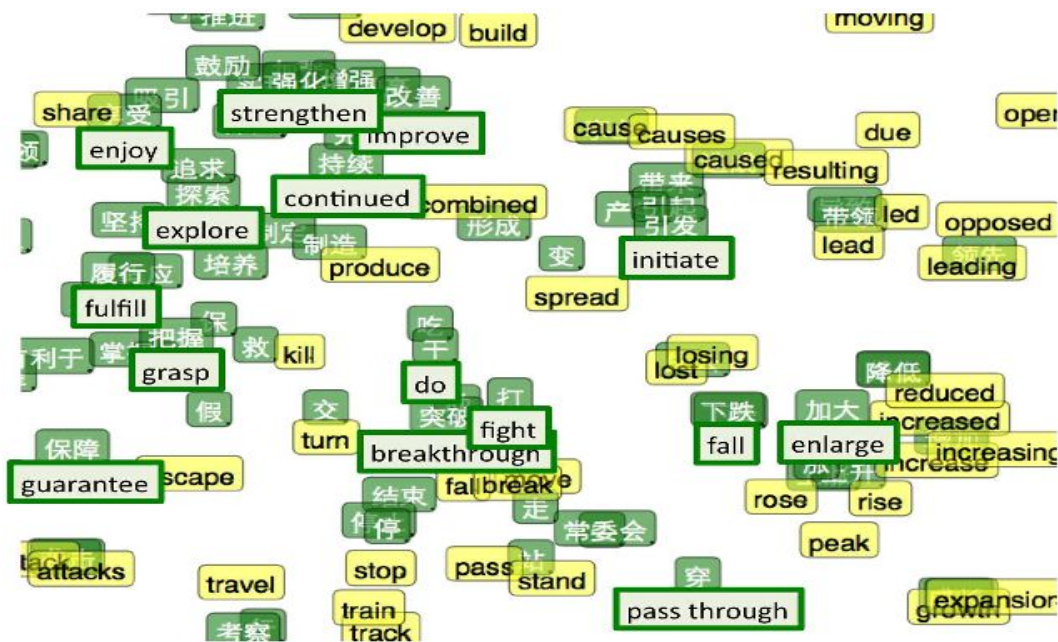
Leverage the inter-connections among languages

English Systems

# Cross-lingual Adaptation

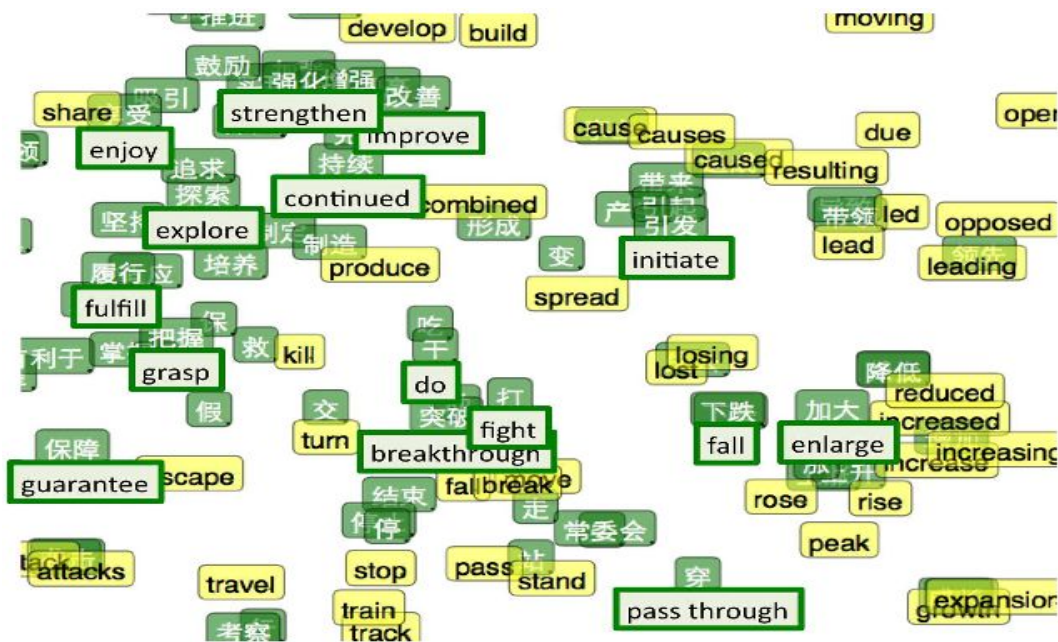
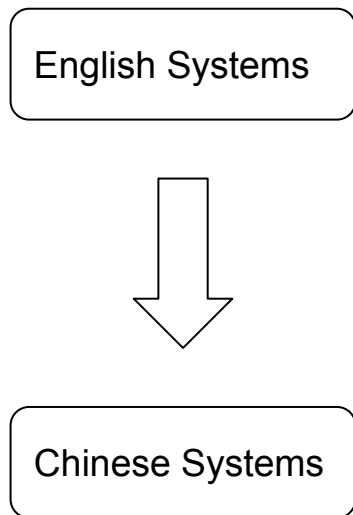
## Leverage the inter-connections among languages

## English Systems



# Cross-lingual Adaptation

## Leverage the inter-connections among languages



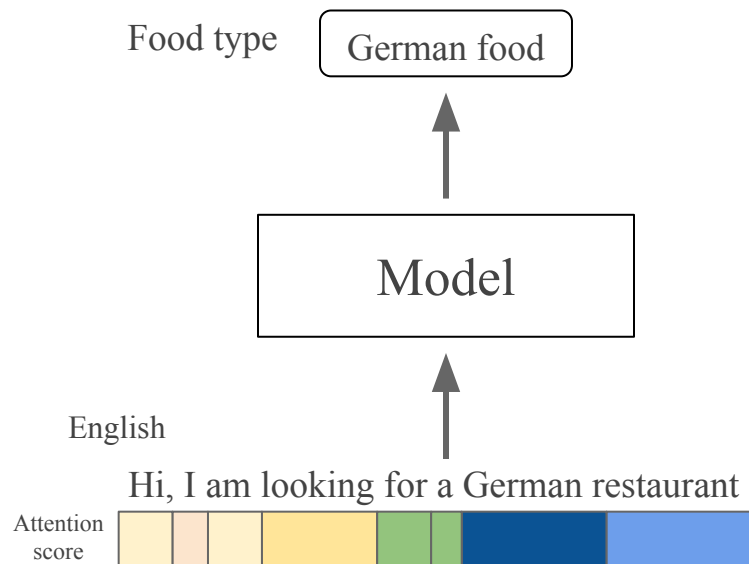
# Related work

- Chen et al (2018)<sup>[1]</sup> utilized large amounts of parallel data or bilingual dictionary to build zero-shot cross-lingual DST systems.
- Schuster et al (2019)<sup>[2]</sup> also leveraged extensive parallel data to build zero-shot cross-lingual NLU systems.
- Collecting bilingual resources is expensive and time-consuming, our work only utilizes very few word pairs as bilingual resources.

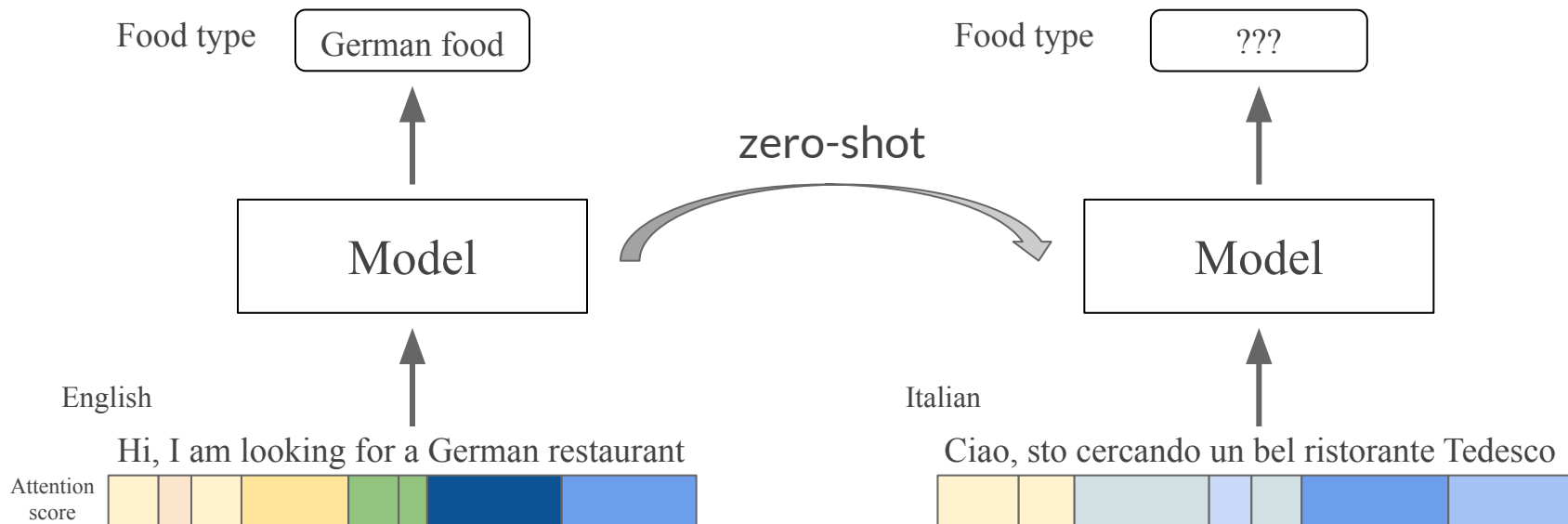
[1] XL-NBT: A Cross-lingual Neural Belief Tracking Framework

[2] Cross-Lingual Transfer Learning for Multilingual Task Oriented Dialog

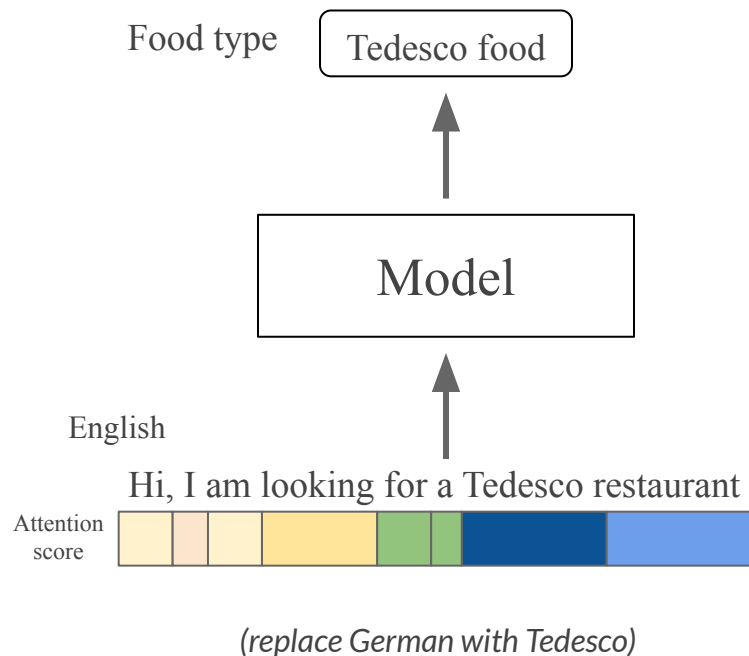
# Intuition of Mixed-Language Training



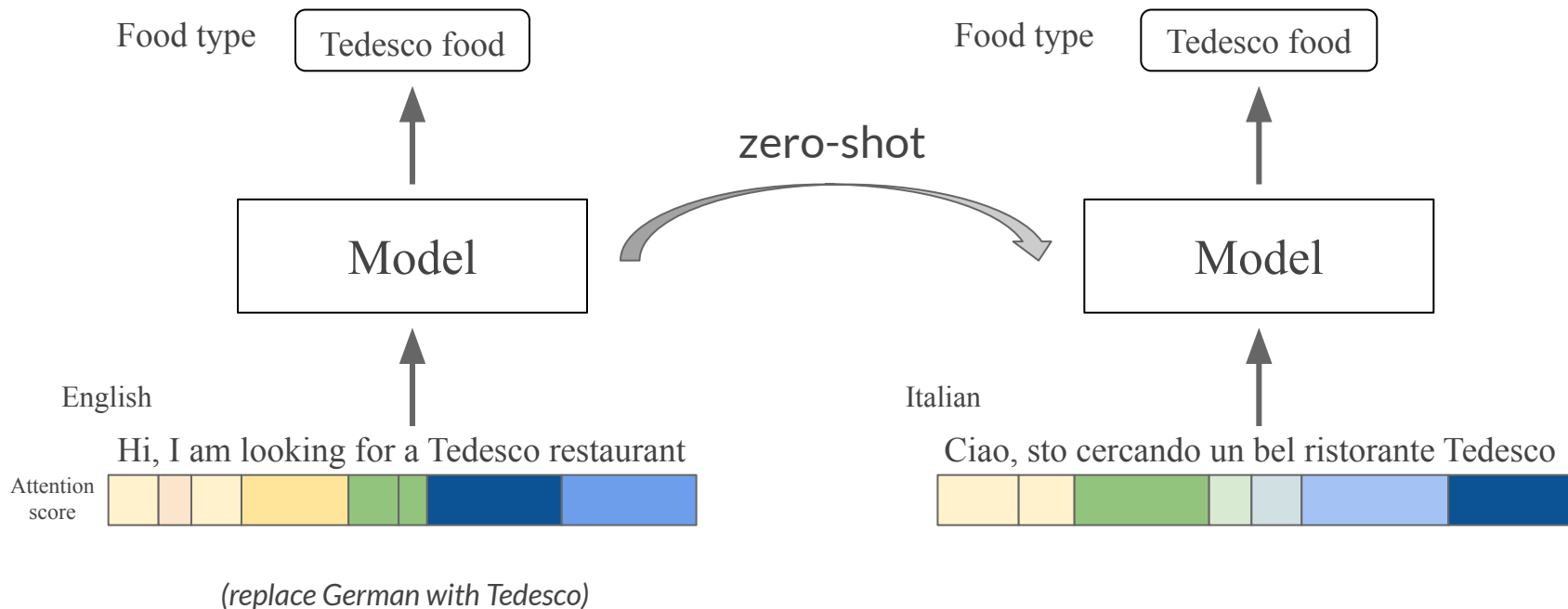
# Intuition of Mixed-Language Training



# Intuition of Mixed-Language Training

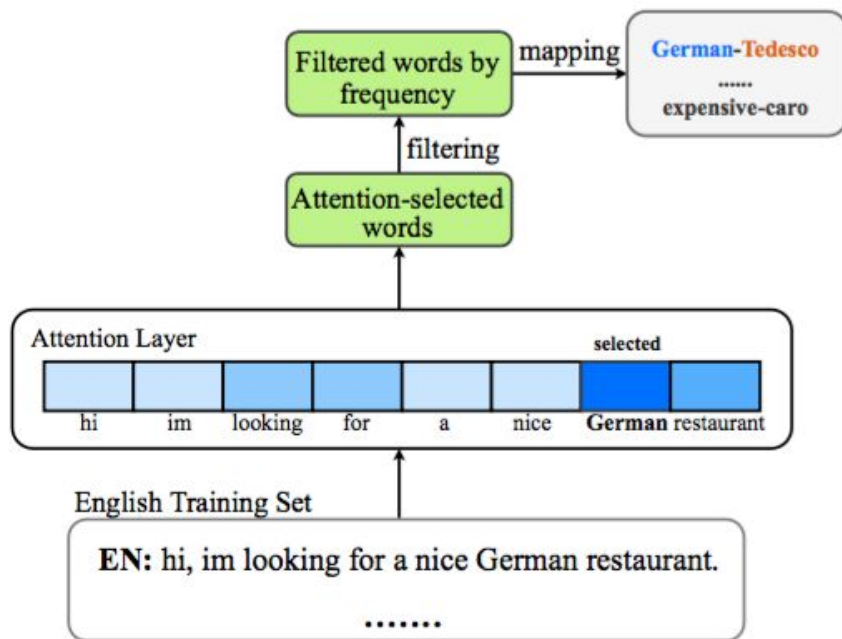


# Intuition of Mixed-Language Training

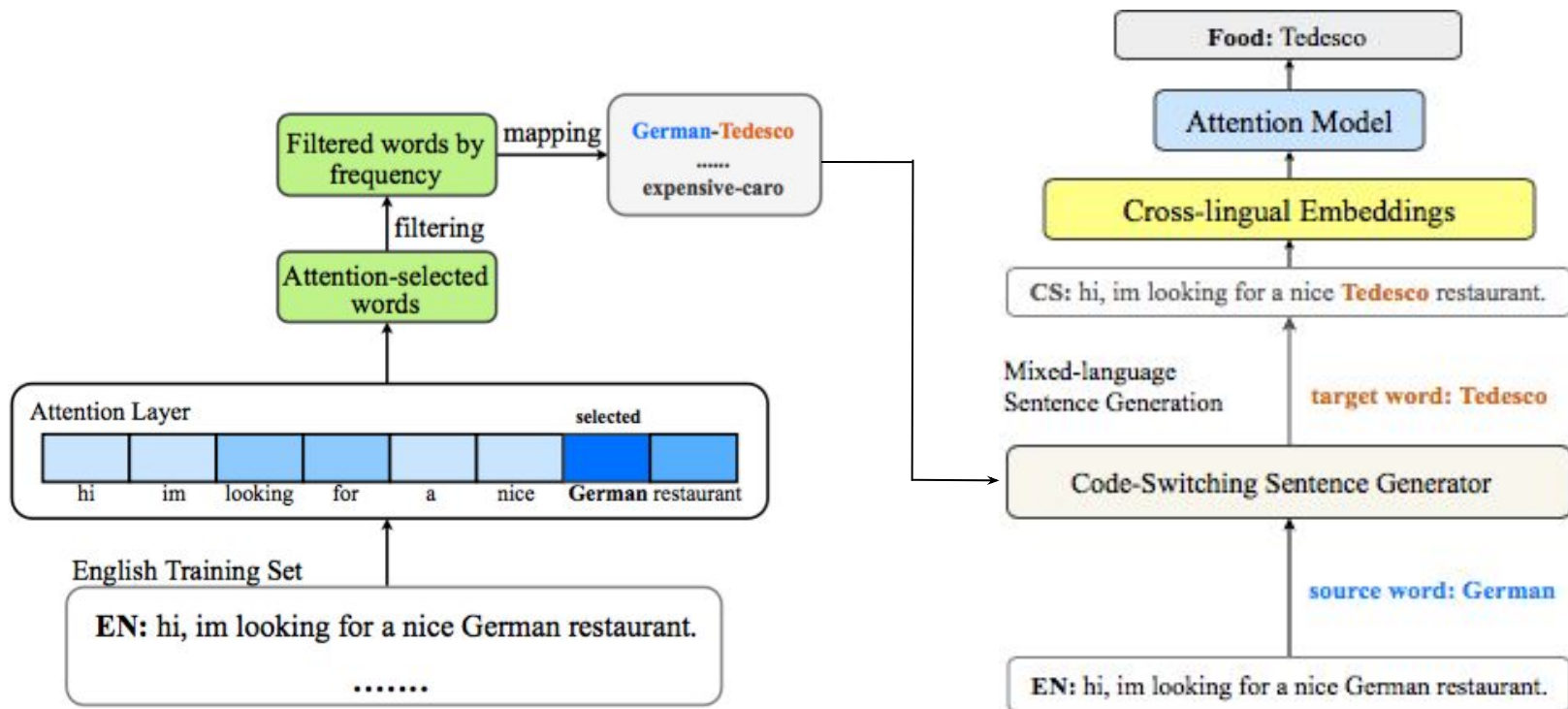




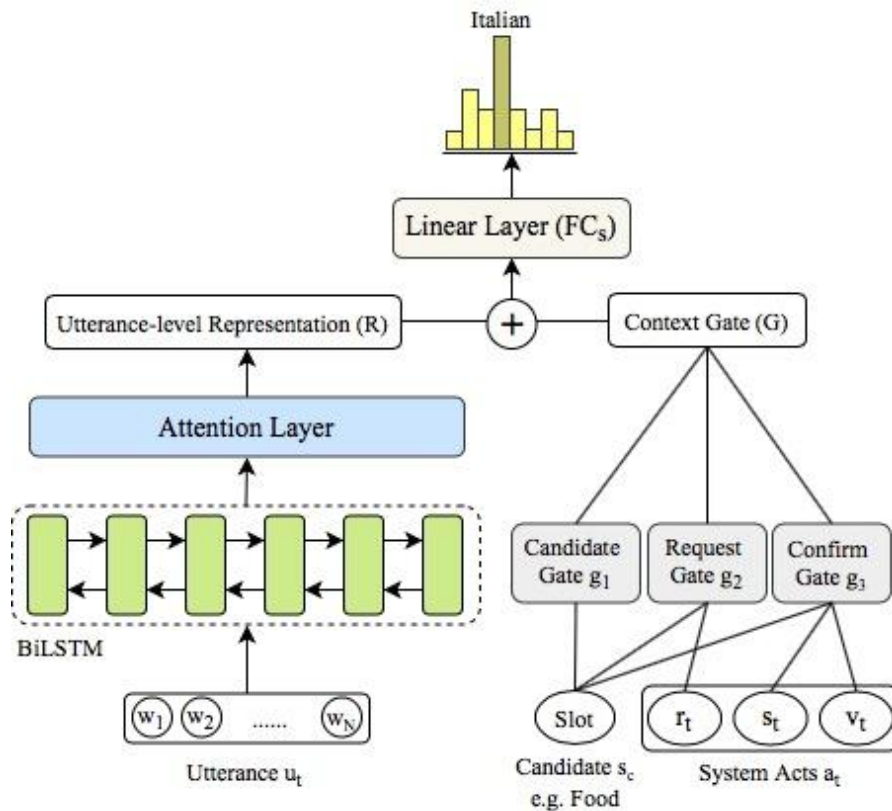
# Attention Layer to select keywords



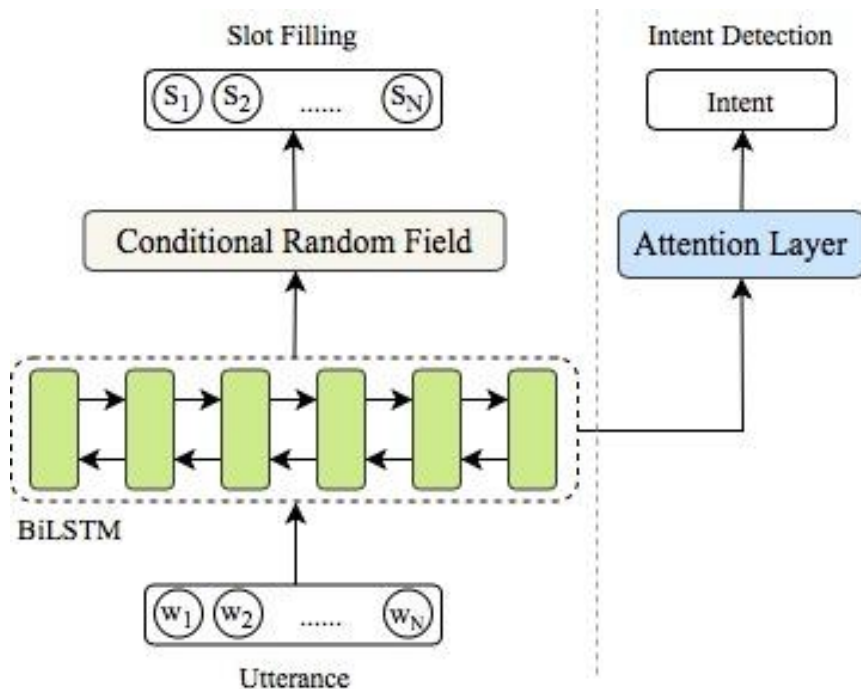
# Mixed-Language Training



# Dialogue State Tracking (DST)



# Natural Language Understanding (NLU)



# Zero-shot Results in DST Task

Model	German								
	slot acc.			joint goal acc.			request acc.		
	BASE	MLT <sub>O</sub>	MLT <sub>A</sub>	BASE	MLT <sub>O</sub>	MLT <sub>A</sub>	BASE	MLT <sub>O</sub>	MLT <sub>A</sub>
MUSE	60.69	68.58	<b>71.38</b>	21.57	30.61	<b>36.51</b>	74.22	80.11	<b>82.99</b>
XLM (MLM)*	52.21	66.26	<b>68.25</b>	14.09	29.45	<b>31.29</b>	75.15	78.48	<b>80.22</b>
+ Transformer	53.81	65.81	<b>68.55</b>	13.97	30.87	<b>32.98</b>	76.83	78.95	<b>81.34</b>
XLM (MLM+TLM)*	58.04	65.39	<b>66.25</b>	16.34	29.22	<b>29.83</b>	75.73	78.86	<b>79.12</b>
+ Transformer	56.52	66.81	<b>68.88</b>	16.59	31.76	<b>33.12</b>	78.56	81.59	<b>82.96</b>
Multi. BERT*	57.61	67.49	<b>69.48</b>	14.95	30.69	<b>32.23</b>	75.31	83.66	<b>86.27</b>
+ Transformer	57.43	68.33	<b>70.77</b>	15.67	31.28	<b>34.36</b>	78.59	84.37	<b>86.97</b>
<i>Ontology Matching<sup>†</sup></i>	24			-			21		
<i>Translate Train<sup>†</sup></i>	41			-			42		
<i>Bilingual Dictionary<sup>‡</sup></i>	51.74			28.07			72.54		
<i>Bilingual Corpus<sup>‡</sup></i>	55			30.84			68.32		
<i>Supervised Training</i>	85.78			78.89			84.02		
	Italian								
Model	slot acc.			joint goal acc.			request acc.		
	BASE	MLT <sub>O</sub>	MLT <sub>A</sub>	BASE	MLT <sub>O</sub>	MLT <sub>A</sub>	BASE	MLT <sub>O</sub>	MLT <sub>A</sub>
MUSE	60.59	73.55	<b>76.88</b>	20.66	36.88	<b>39.35</b>	79.09	82.24	<b>84.23</b>
Multi. BERT*	53.34	65.49	<b>69.48</b>	12.88	26.45	<b>31.41</b>	76.12	84.58	<b>85.18</b>
+ Transformer	54.56	66.87	<b>71.45</b>	12.63	28.59	<b>33.35</b>	77.34	82.93	<b>84.96</b>
<i>Ontology Matching<sup>†</sup></i>	23			-			21		
<i>Translate Train<sup>†</sup></i>	48			-			51		
<i>Bilingual Dictionary<sup>‡</sup></i>	73			39.01			77.09		
<i>Bilingual Corpus<sup>‡</sup></i>	72			41.23			81.23		
<i>Supervised Training</i>	88.92			80.22			91.05		

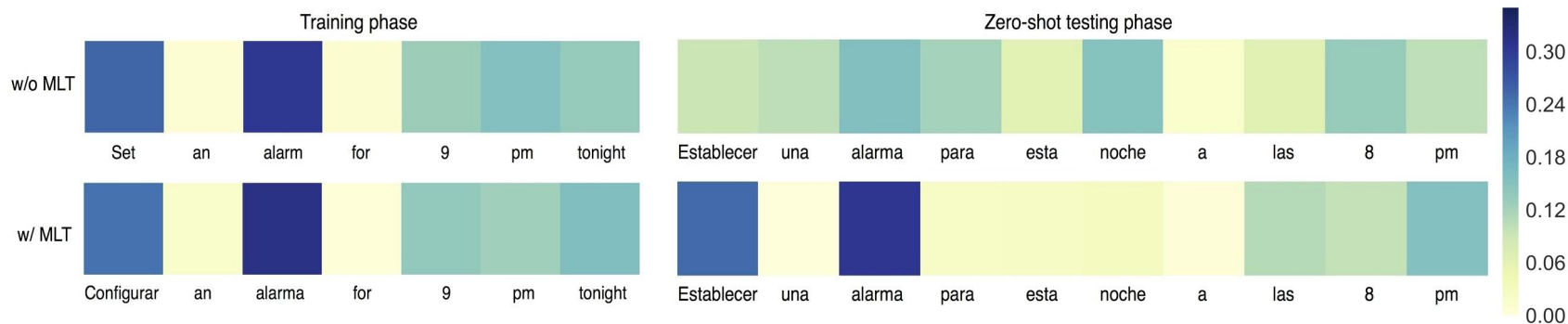
Zero-shot results for the target languages on Multilingual WOZ 2.0. MLT<sub>A</sub> denotes our approach (attention-informed MLT), which utilizes the same number of word pairs (90 word pairs) as MLT<sub>O</sub> (MLT based on ontology).

# Zero-shot Results in NLU Task

Model	Spanish						Thai					
	Intent acc.			Slot F1			Intent acc.			Slot F1		
	BASE	MLT <sub>H</sub>	MLT <sub>A</sub>	BASE	MLT <sub>H</sub>	MLT <sub>A</sub>	BASE	MLT <sub>H</sub>	MLT <sub>A</sub>	BASE	MLT <sub>H</sub>	MLT <sub>A</sub>
RCSLS	37.67	77.59	<b>87.05</b>	22.23	<b>59.12</b>	57.75	35.12	68.63	<b>81.44</b>	8.72	29.44	<b>30.42</b>
XLM (MLM)	60.8	75.11	<b>83.95</b>	38.55	63.29	<b>66.11</b>	37.59	46.34	<b>65.31</b>	8.12	19.03	<b>20.43</b>
+ Transformer	62.33	82.83	<b>85.63</b>	41.67	66.53	<b>67.95</b>	40.31	57.27	<b>68.55</b>	11.45	26.02	<b>27.45</b>
XLM (TLM+MLM)	62.48	81.34	<b>84.91</b>	42.27	65.71	<b>66.48</b>	31.62	50.34	<b>65.25</b>	7.91	19.22	<b>19.88</b>
+ Transformer	65.32	83.79	<b>87.48</b>	44.39	66.03	<b>68.55</b>	37.53	68.62	<b>72.59</b>	12.84	26.56	<b>27.98</b>
Multi. BERT	73.73	77.51	<b>86.54</b>	51.73	<b>74.51</b>	74.43	28.15	52.25	<b>70.57</b>	10.62	24.41	<b>28.47</b>
+ Transformer	74.15	82.9	<b>87.88</b>	54.28	<b>74.88</b>	73.89	26.54	53.84	<b>73.46</b>	11.34	26.05	<b>27.12</b>
<i>Zero-shot SLU<sup>†</sup></i>		46.64			15.41			35.64			12.11	
<i>Multi. CoVe</i>		53.34			22.50			66.35			32.52	
<i>Multi. CoVe w/ auto</i>		53.89			19.25			70.70			35.62	
<i>Translate Train</i>		85.39			72.87			95.85			55.43	

Zero-shot results on multilingual NLU dataset (Schuster et al. 2019), and the number of word pairs on both MLT<sub>H</sub> and MLT<sub>A</sub> is 20.

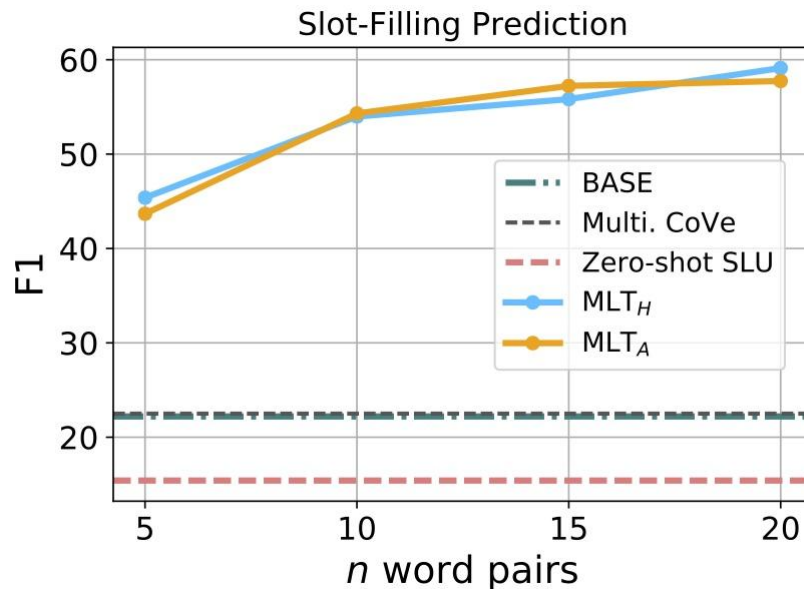
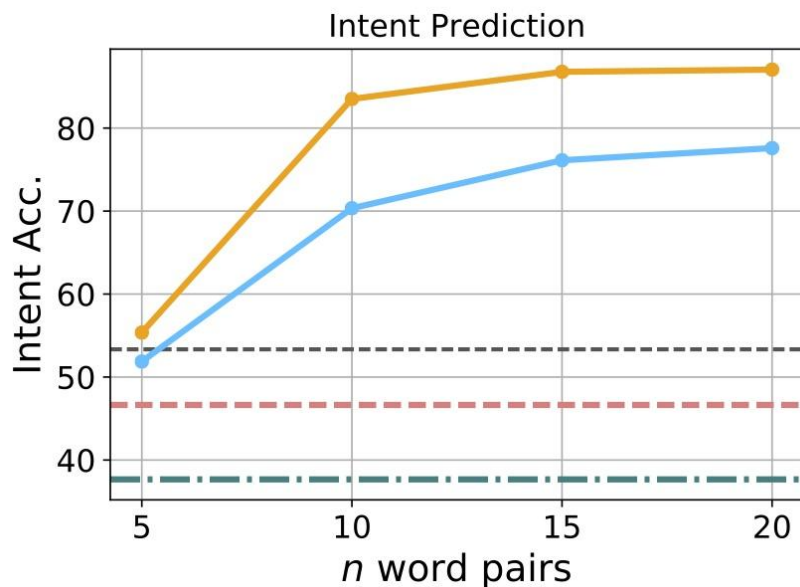
# Visualization



Attentions on words in both training and testing phases.



# Zero-shot Results in NLU Task



The dynamics of the NLU task: intent and slot-filling results with different numbers of word pairs on Spanish test data.



# Conclusion

- We propose attention-informed mixed-language training for cross-lingual task-oriented dialogue systems.
- Our approach utilizes very few task-related parallel word pairs base on the attention scores.
- The task-related words have a generalization ability to other words that have similar semantics in target languages.

# Thanks!



**Check our code**

<https://github.com/zliucr/mixed-language-training>