

Towards multilingual sequence labelling using cross-lingual word embeddings

*Chao Li, Carsten van Weelden, Luigi Lorato, Carsten L. Hansen,
Mihai Rotaru and Jakub Zavrel*

Today

- Textkernel
- BI-LSTM Sequence labeling Model
- Embedding Alignment and Multilingual model
- Conclusion and Future work

Textkernel

- Focus: NLP applications in HR domain
 - CV Parsing for 17 languages
 - Job parsing for 7 languages
 - Semantic search
 - Recommended jobs and candidates
- Started in 2001, located in Amsterdam
- 100+ people, +30 nationalities

Outline

- Textkernel
- **BI-LSTM Model for Sequence Labeling**
- Embedding Alignment and Multilingual model
- Conclusion and Future work

CV

Fiona Lee City: Paris
Phone: 0965.....

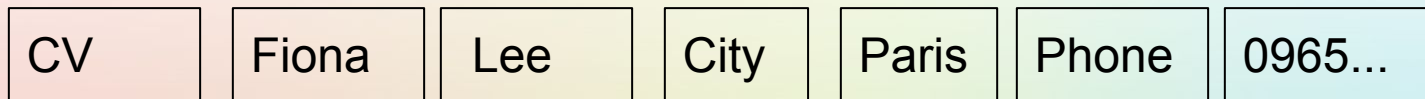
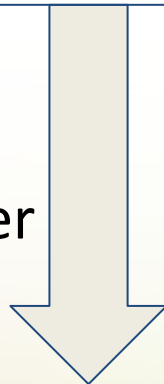
Work Experience

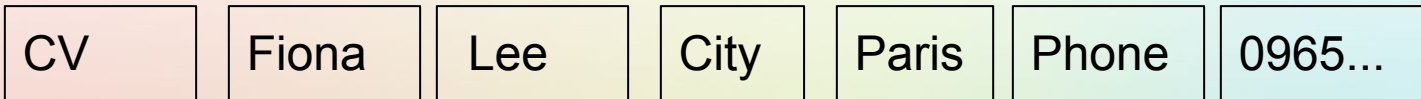
1996-2000

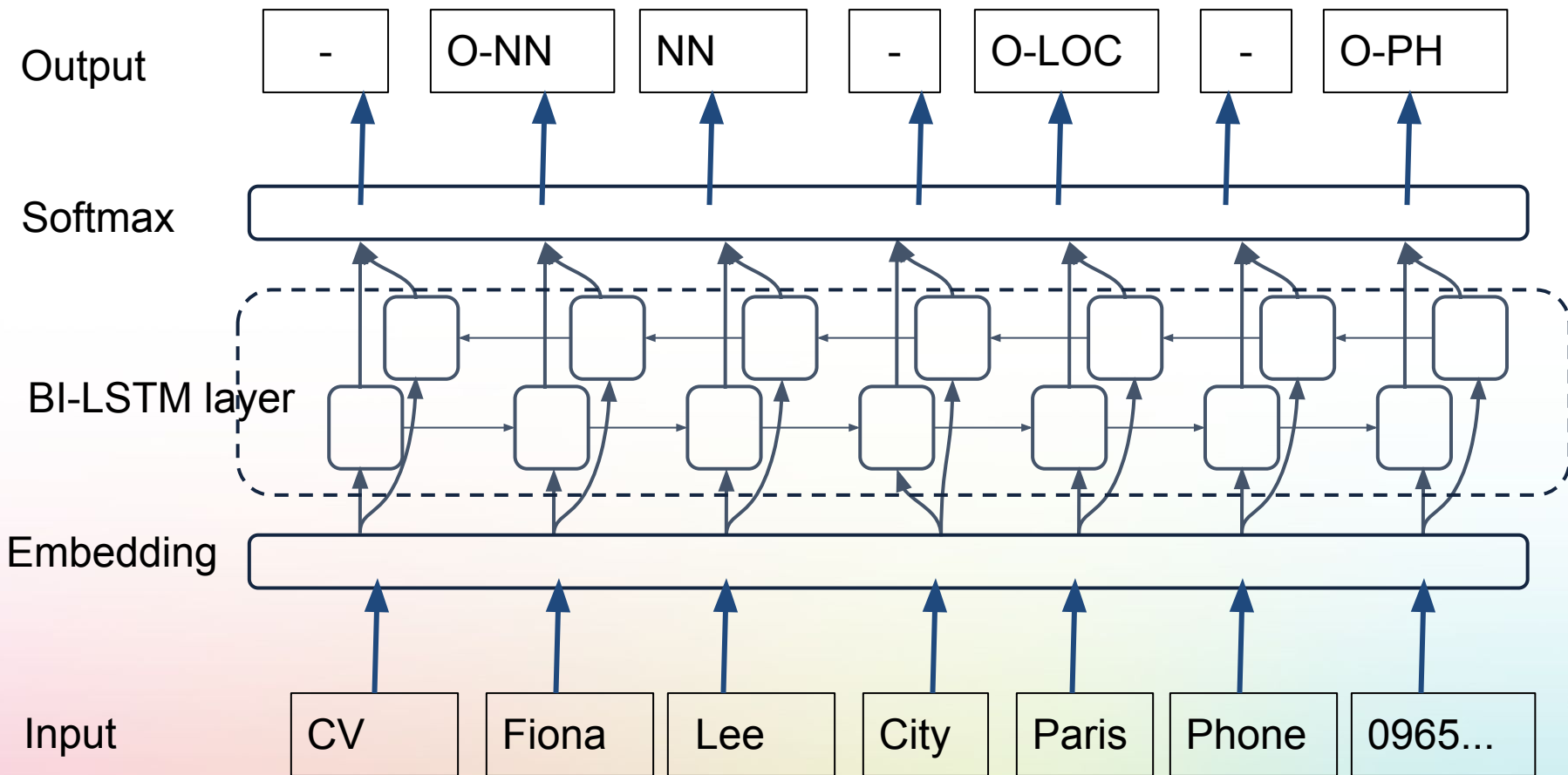
Java Engineer

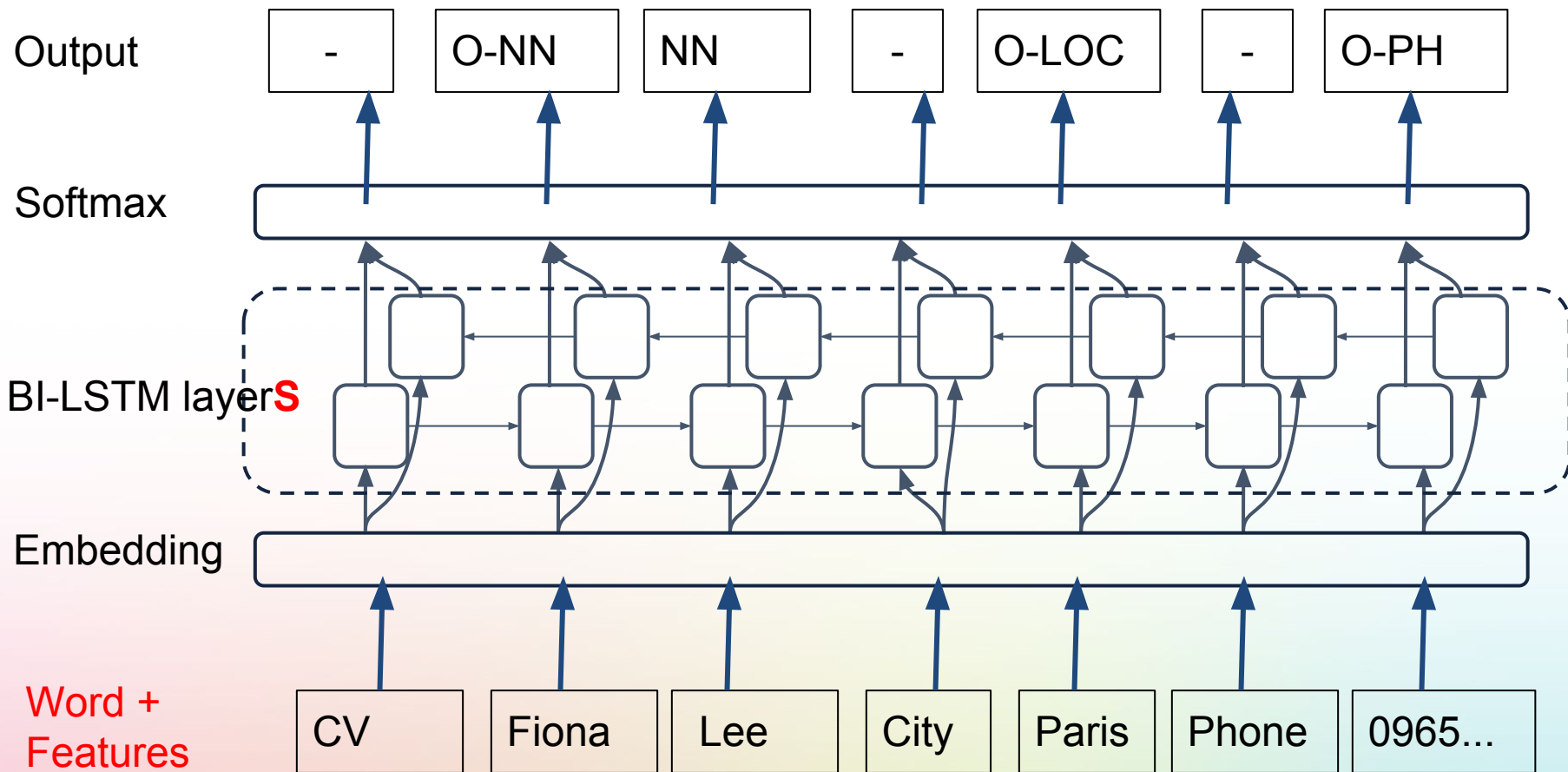
IBM, Amsterdam

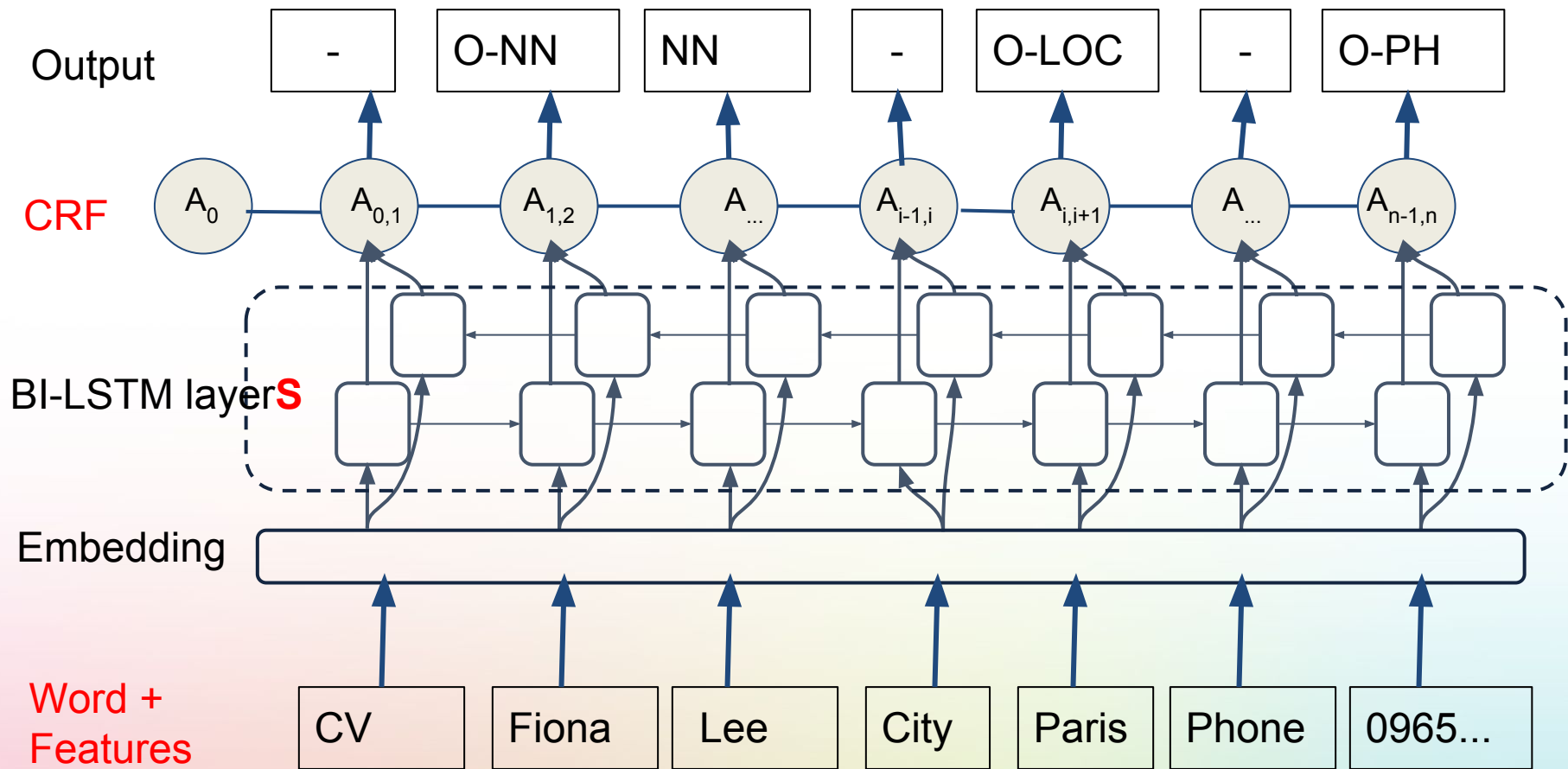
.....

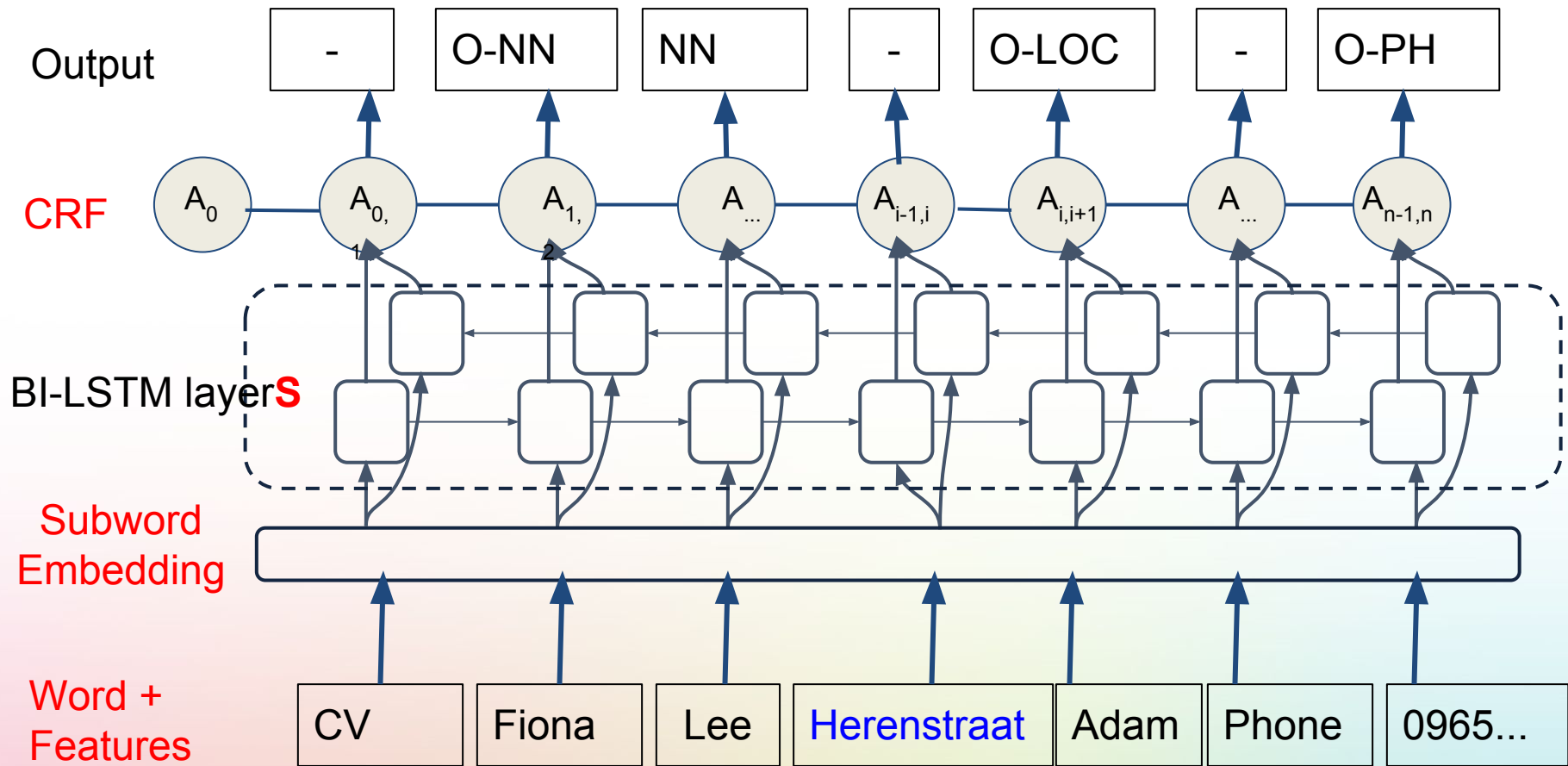






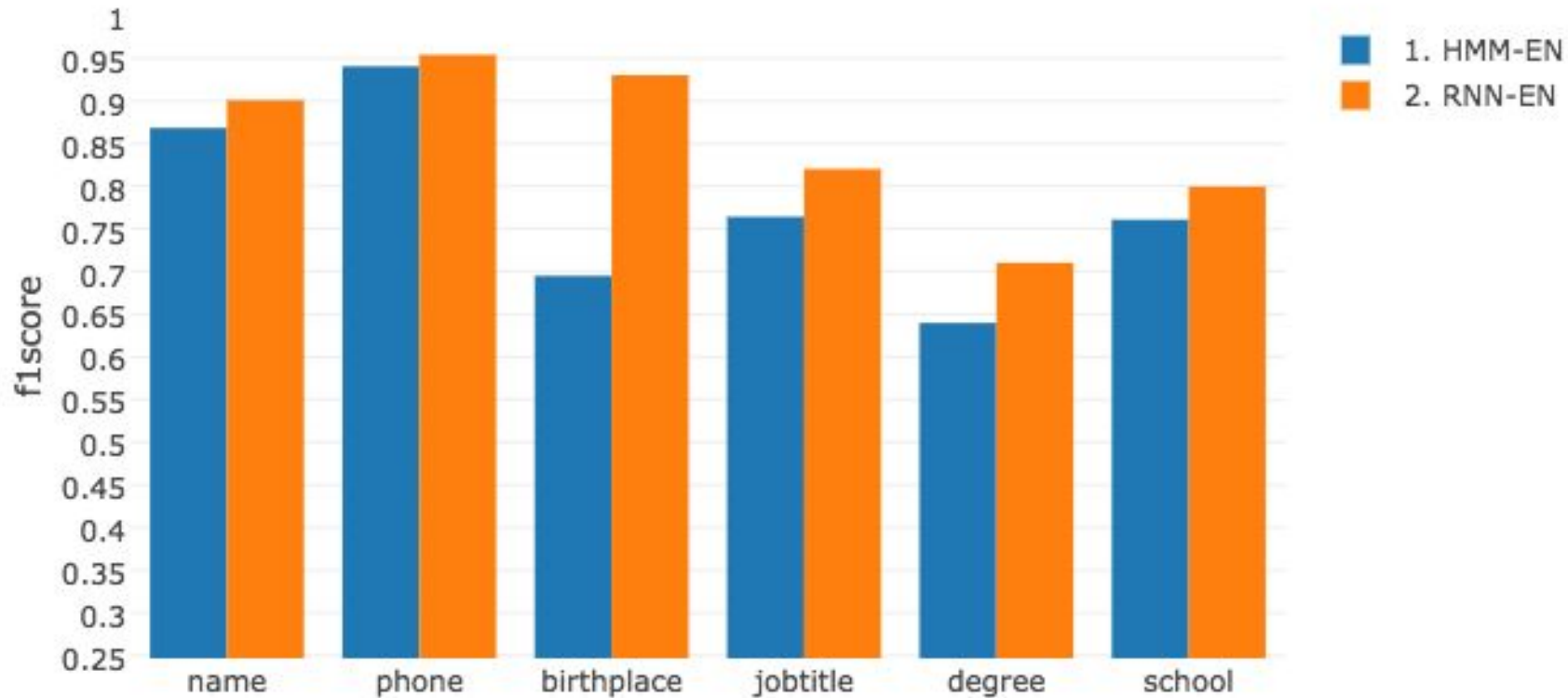






NN Sequence labeling Model

- word based
 - Multilayer BI-LSTM (dropout, batchNorm)
 - CRF layer to improve the label consistency
 - Subword embedding to eliminate unknown words
- Phrase based
 - CNN to compose phrase embeddings



Outline

- Textkernel
- BI-LSTM Model for Sequence Labeling
- **Embedding Alignment and Multilingual model**
- Conclusion and Future work

Multilingual CV parsing model

- Hard to get training data for new languages
- CVs written similarly across languages/countries

简历

菲娜 李

城市: 巴黎

电话: 0965.....

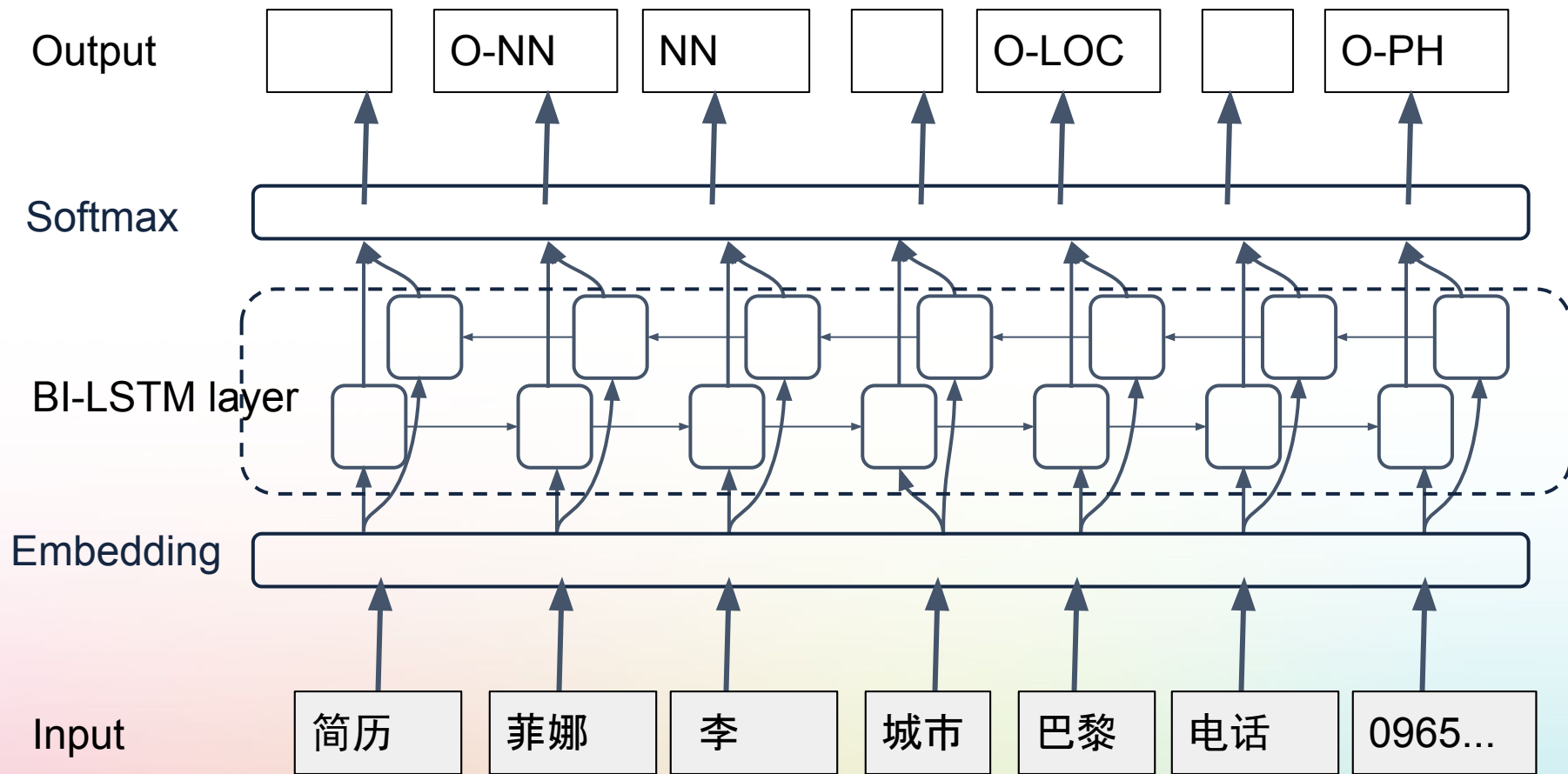
工作经历

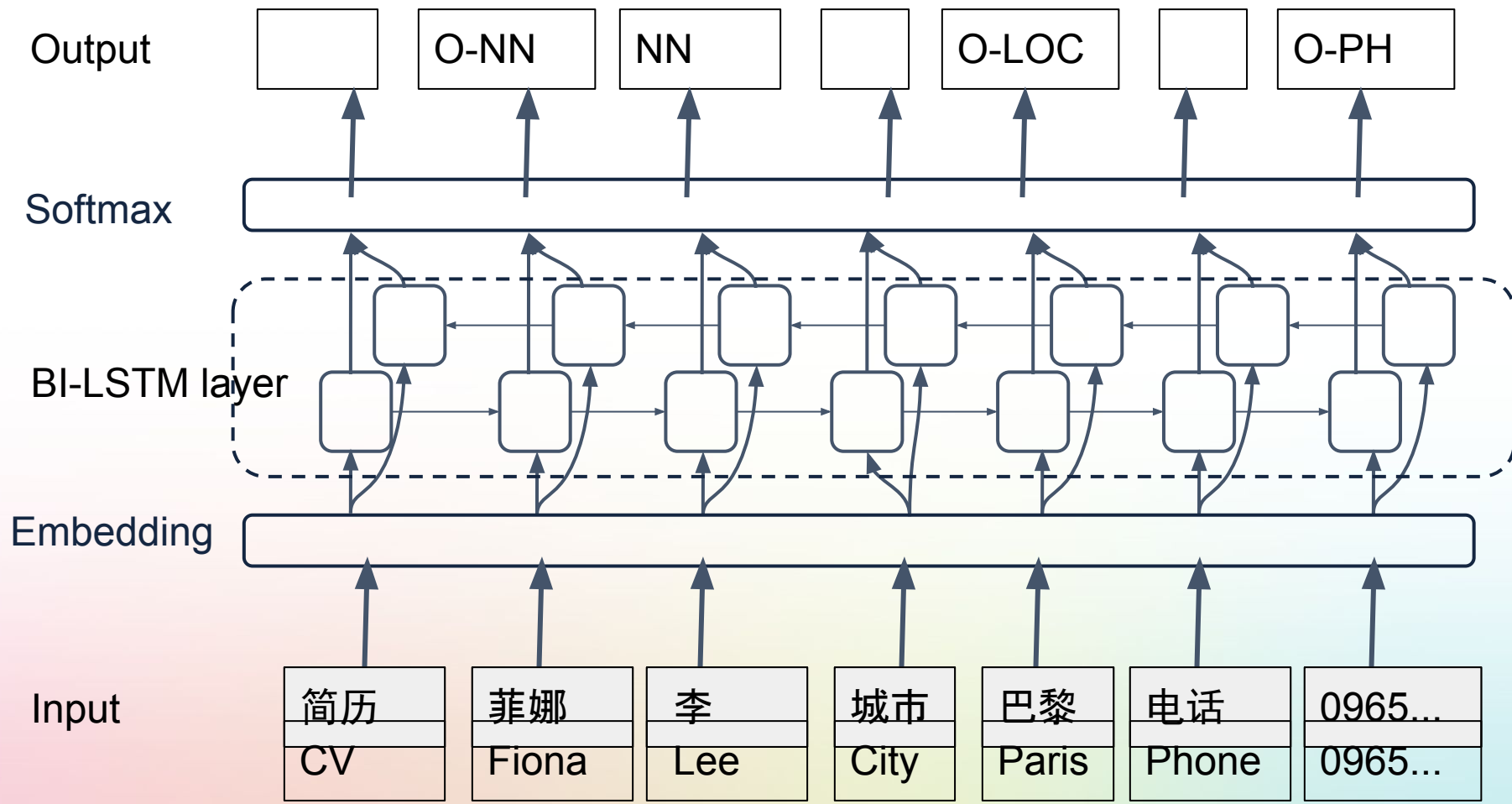
1996-2000

Java 工程师

IBM, 阿姆斯特丹

.....





Multilingual CV parsing model

- CVs written similarly across languages/countries
- Align embeddings across languages:

CN to EN: Find \mathbf{U} and \mathbf{V} such that S_{en} and S_{cn} are maximum correlated:

- $S_{en} = \mathbf{U} * \mathbf{WV}_{en}$
- $S_{cn} = \mathbf{V} * \mathbf{WV}_{cn}$

Multilingual CV parsing model

- CVs written similarly across languages/countries
- Align embeddings across languages

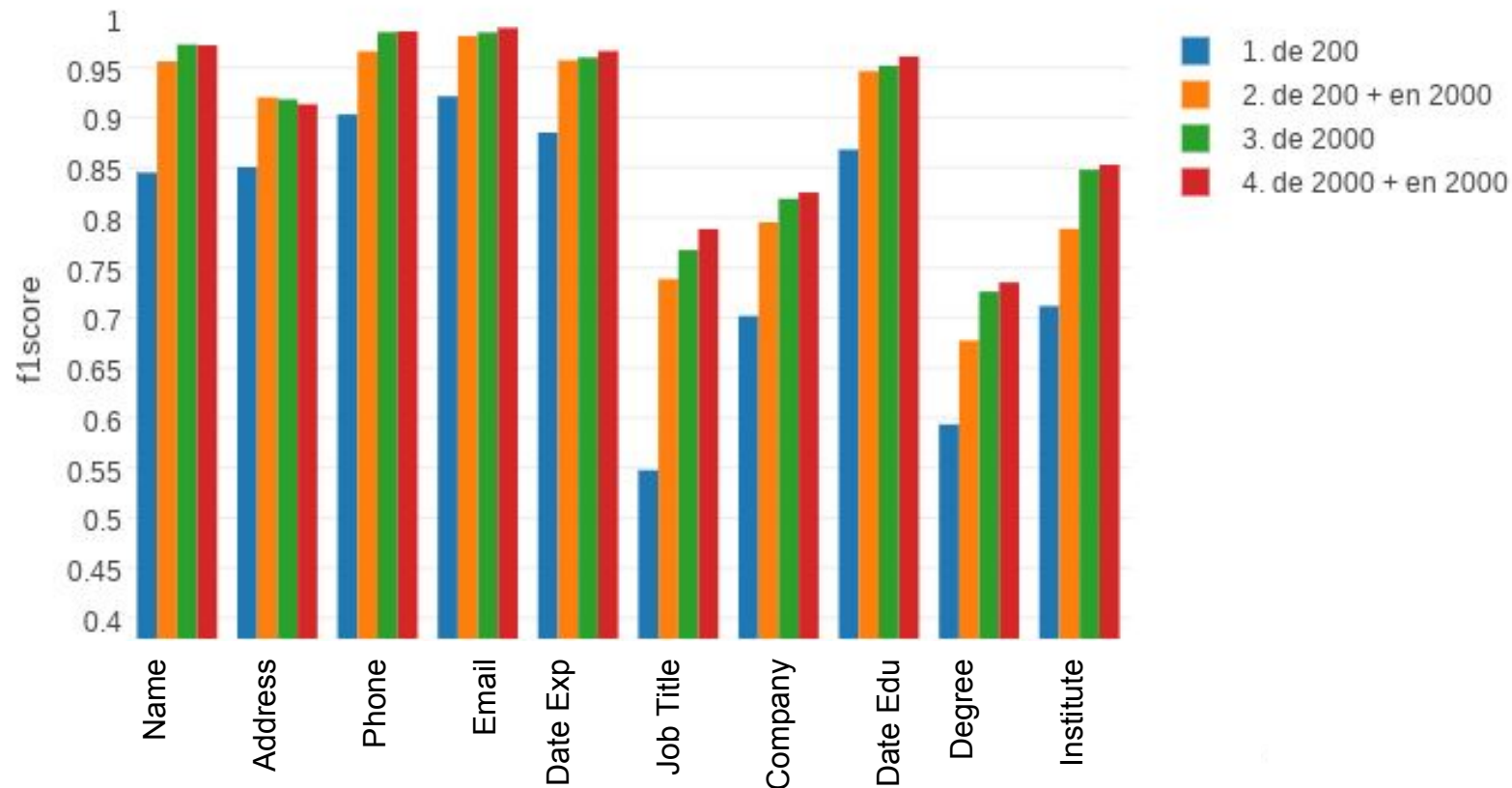
CN to EN: mapping Chinese embedding to English vector space

- $S_{en}^* = U^{-1} * U * WV_{en} = WV_{en}$
- $S_{cn}^* = U^{-1} * V * WV_{cn}$

Multilingual CV parsing model

- CVs written similarly across languages/countries
- Align embeddings across languages

[Embedding Alignment](#)



Summary

Conclusion:

- RNN sequence labeling model is better than baseline models
- Aligned embedding allow cross lingual training
- Pave the way for a multilingual model

Outlook

- Improve the alignment of embeddings
- Better evaluation of multilingual embeddings
- Train a multilingual model and test on an unseen language

Thank you?