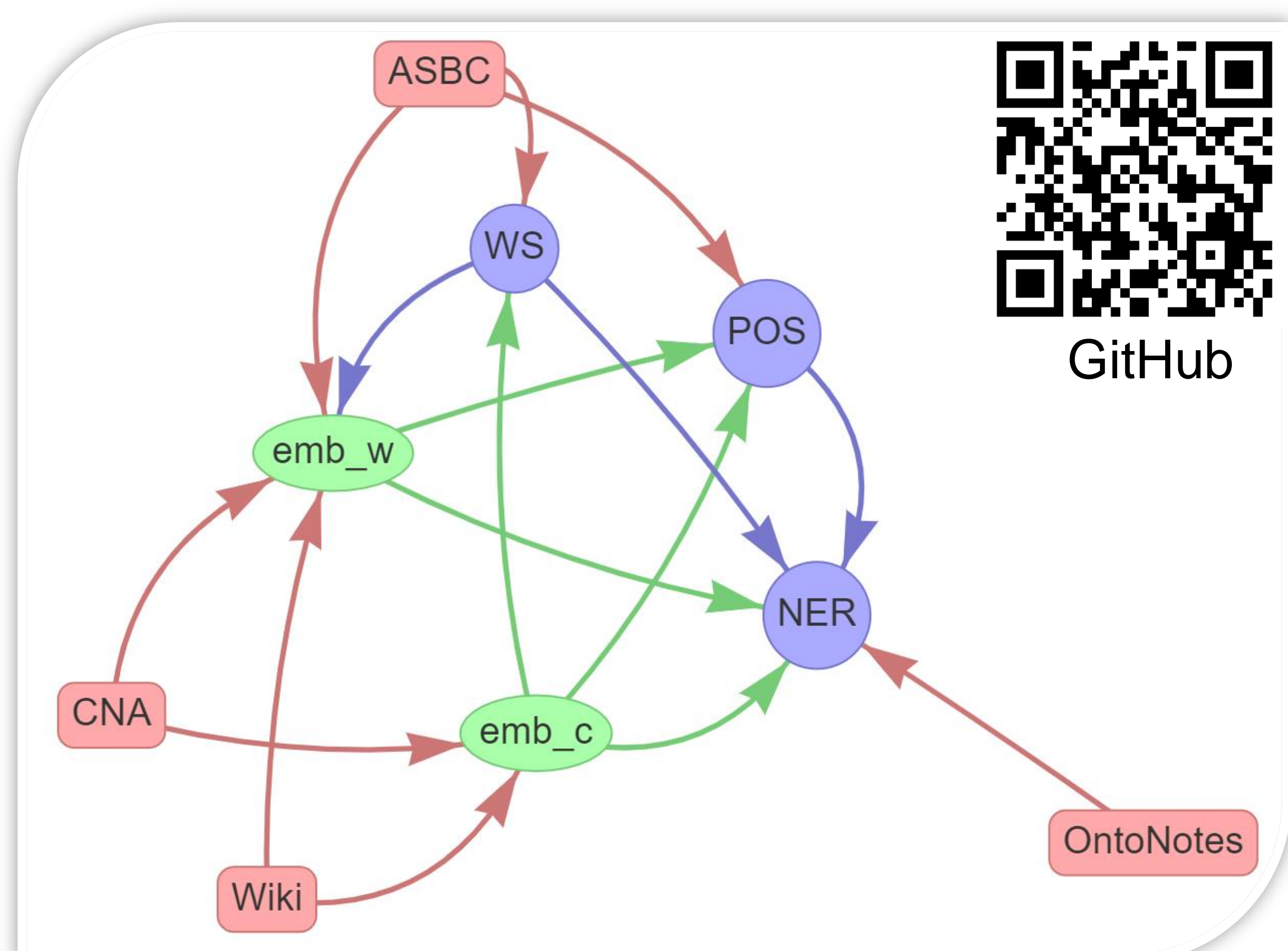


建構概念為本且具語義結合性的中文知識庫

馬偉雲助研究員團隊 中研院資訊所

CkipTagger: 開源的斷詞(Ws)、詞性標注(POS)、實體辨識(NER)工具



18類實體辨識

(人,事件,日期,地理,產品,組織,法律,時間,行政區,設施,社會團體,語言,藝術品,錢,數字,序數,數量,百分比)

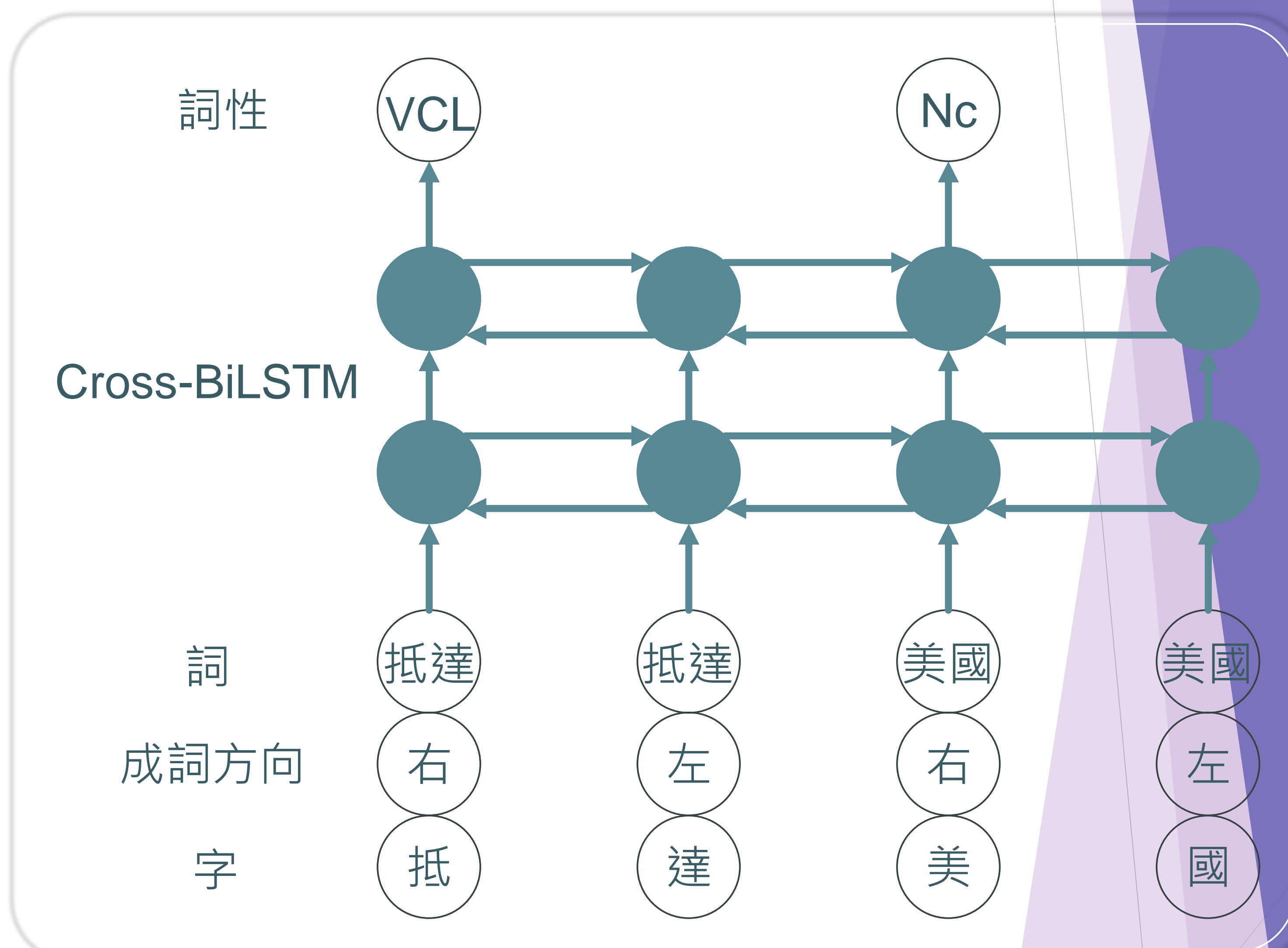
傅達仁_{PERSON}今將執行安樂死，卻突然爆出自己20年前_{DATE}遭緯來體育台_{ORG}封殺，他不懂自己哪裡得罪到電視台。

61類詞性標注

(各類名詞,各類動詞,各類副詞,各類標點,...)

傅達仁_(Nb) 今_(Nd) 將_(D) 執行_(VC) 安樂死_(Na) , _(COMMACATEGORY) 卻_(D) 突然_(D) 爆出_(VJ) 自己_(Nh) 20_(Neu) 年_(Nf) 前_(Ng) 遭_(P) 緯來_(Nb) 體育台_(Na) 封殺_(VC) , _(COMMACATEGORY) 他_(Nh) 不_(D) 懂_(VK) 自己_(Nh) 哪裡_(Ncd) 得罪到_(VJ) 電視台_(Nc) ° _(PERIODCATEGORY)

詞性標注模型



語料

Corpus		Size	Training
CNA	中央社新聞	11億字	字向量/詞向量
Wiki	中文維基	4.6億字	字向量/詞向量
ASBC	中研院平衡語料庫	1633萬字	斷詞/詞性標注
OntoNotes	多重標記語料庫	151萬字	實體辨識

實驗分析

	ASBC-test			OntoNotes-Test	
	WS	POS	NER		
	Prec.	Rec.	F1	Acc.	F1
CkipTagger POS	--	--	--	97.20	--
CkipTagger WS+POS+NER	97.49	97.17	97.33	94.59	77.98
CKIPWS	95.85	95.96	95.91	90.62	--
Jieba_zhTW	90.51	89.10	89.80	--	--

$$*POS\ acc = \frac{\# predicted\ words\ with\ correct\ boundary\ and\ POS\ tags}{\# gold\ standard\ words}$$

論文

Peng-Hsuan Li, Tsu-Jui Fu, and Wei-Yun Ma. 2019. Remediating BiLSTM-CNN Deficiency in Modeling Cross-Context for NER. arXiv preprint arXiv:1908.11046.