

Terminologies augmented recurrent neural network for clinical named entity recognition

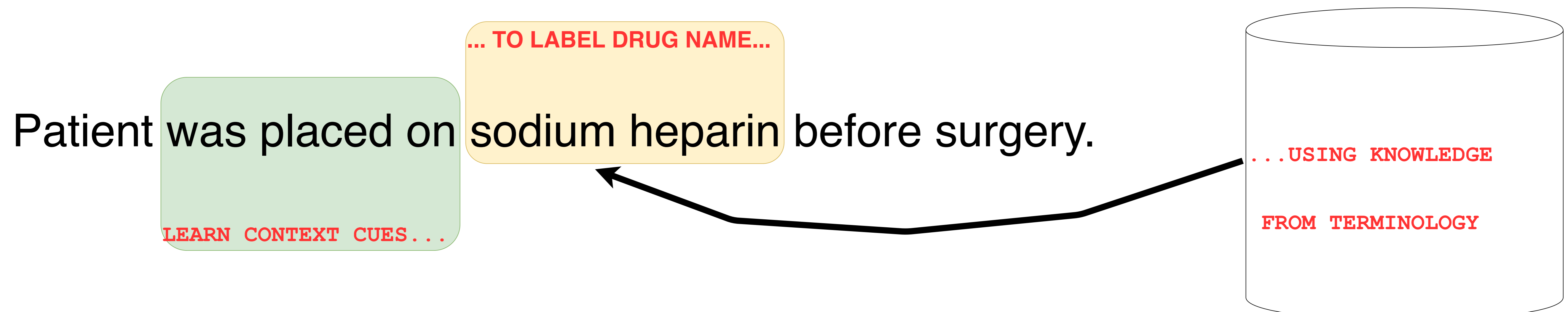
Ivan Lerner, MSc¹⁻²; Nicolas Paris, MSc²⁻³; Xavier Tannier, PhD⁴

¹ INSERM, Centre de Recherche des Cordeliers, UMR 1138 Equipe 22, Paris Descartes, Sorbonne Paris Cité University, Paris, France

² AP-HP, DSI-WIND, Paris, France

³ LIMSI, CNRS, Univ. Paris-Sud, Université Paris-Saclay, F-91405, Orsay, France

⁴ Sorbonne Université, Inserm, Univ Paris 13, LIMICS, F-93017 Bobigny, France

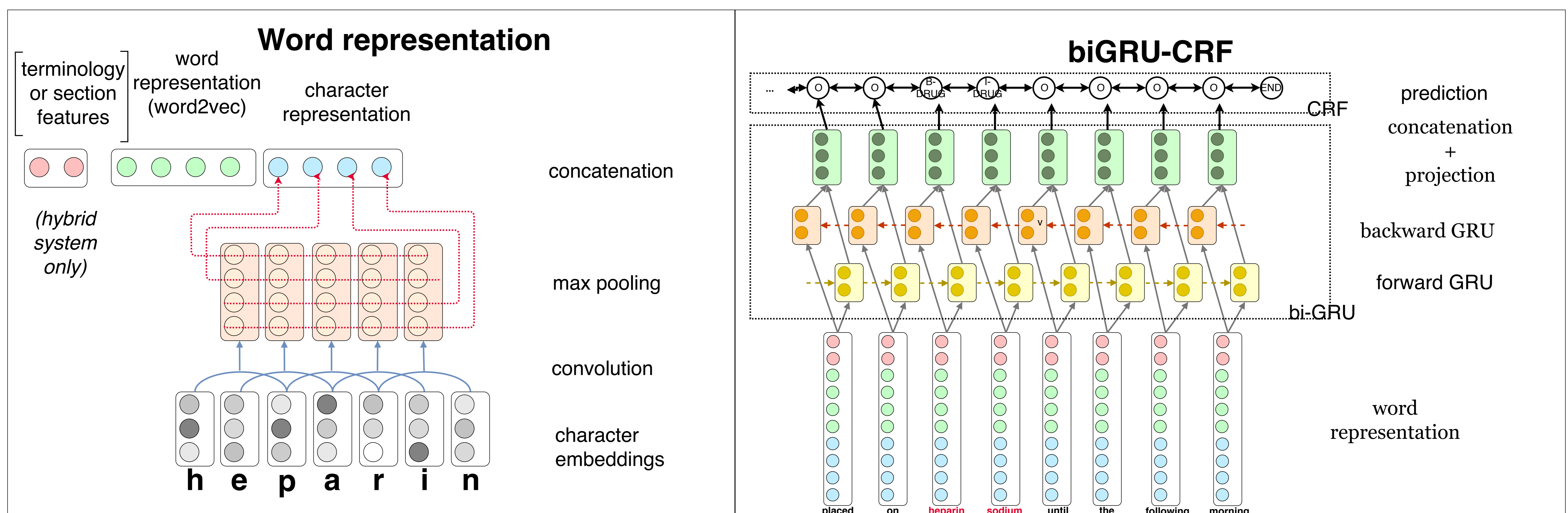


Datasets

- english: i2b2-2009, 268 documents, 8573 drug names
- french: APcNER, 147 documents, 1076 drug names

Terminologies:

- French: ATC, BPDM, CCAM, CIM-10, DRC, UMLS® (including SNOMED 3.5 CT®), GLIMS, QDOC
- English: UMLS® (including SNOMED 3.5 CT®)



Results

Inter-annotator agreement (10 documents)			APcNER corpus description (147 documents)			
Entity types	Exact F-measure	Non exact F-measure	# entities	n-grams (%)		
				n = 1	n = 2	n ≥ 3
Drug name	.85	.92	1076	1014 (94)	54 (.5)	8 (.1)
Sign or symptom	.55	.71	432	356 (82)	65 (15)	11 (.3)
Disease or disorder	.65	.77	1672	1238 (74)	330 (20)	104 (.6)
Diagnostic procedure or lab test	.70	.87	1156	808 (70)	297 (27)	51 (.4)
Therapeutic procedure	.51	.71	501	414 (83)	73 (15)	14 (.3)

Corpus	System	Exact-match			Partial-match		
		F ^a [min-max]	P ^b [min-max]	R ^c [min-max]	F ^a [min-max]	P ^b [min-max]	R ^c [min-max]
i2b2-2009	Terminologies	73.0	76.7	69.7	84.6	88.9	80.6
	biGRU-CRF	91.1 [90.0-91.8]	90.6 [89.7-92.6]	91.7 [87.6-93.1]	93.5 [92.4-94.1]	92.9 [92.2-95.0]	94.2 [90.0-96.0]
	Hybrid system	92.2 [91.2-93.0]	92.1 [91.4-93.1]	92.2 [90.5-93.8]	94.7 [94.3-95.2]	94.6 [94.0-95.2]	94.7 [93.5-96.1]
APcNER	Terminologies	75.0	70.8	79.7	77.7	73.3	82.5
	biGRU-CRF	81.9 [81.2-82.4]	86.6 [84.9-88.7]	77.8 [76.6-78.9]	86.4 [85.1-87.7]	91.4 [90.3-93.4]	82.0 [80.1-84.0]
	Hybrid system	86.4 [86.2-86.8]	89.6 [87.7-90.9]	83.4 [82.3-84.7]	90.4 [89.9-91.1]	93.8 [91.9-94.9]	87.2 [85.9-88.6]

^aF-measure; ^bPrecision; ^cRecall

Conclusion



- APcNER: 28 hours of annotation for 5 entities for 147 documents
- APcNER annotation allowed 82% F1 for drug name
- Terminology features boost model performance (+4.3% F1)
- Terminology features boost is smaller (+1%) with more data (268 documents)

