
Effective Feature Representation for Clinical Text Concept Extraction

Yifeng Tao^{1,2}, Bruno Godefroy¹, Guillaume Genthial¹, Christopher Potts^{1,3,*}

¹Roam Analytics

²Carnegie Mellon University

³Stanford University

ROAM



Carnegie Mellon University
School of Computer Science



Stanford
University

Background: Healthcare Text Datasets

- Crucial information of healthcare recorded only in free-form text

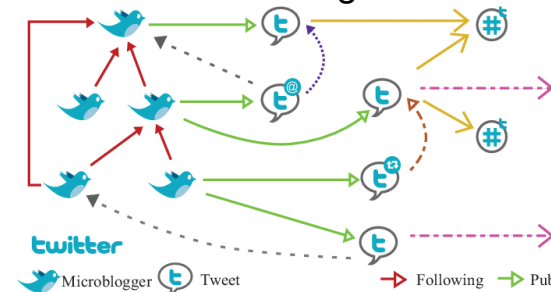
Scientific
Chemical-Disease Relations



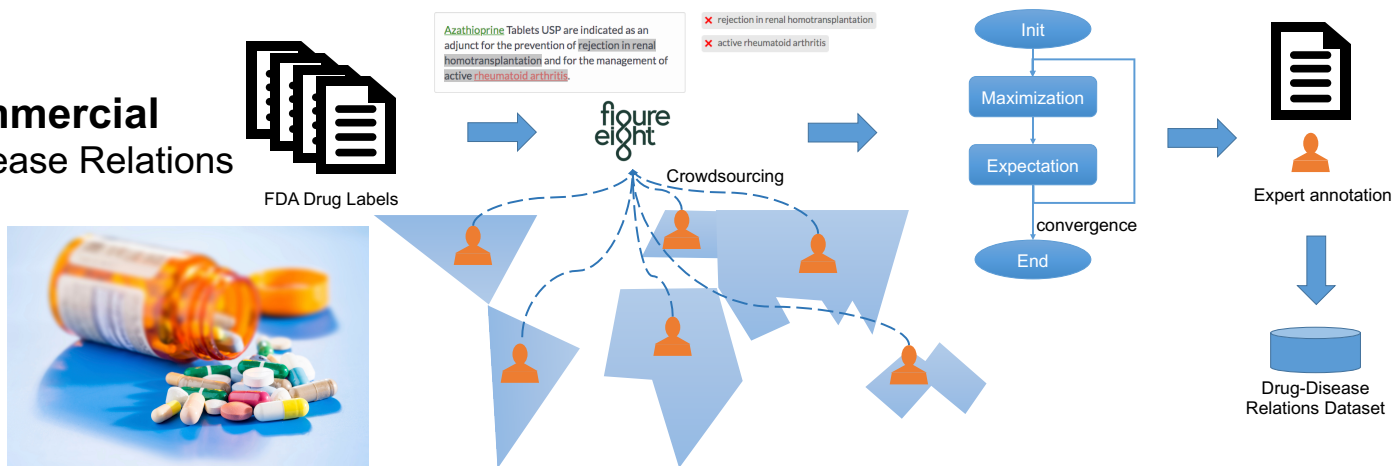
Clinical
Diagnosis Detection
Prescription Reasons



Social
Penn Adverse Drug Reactions



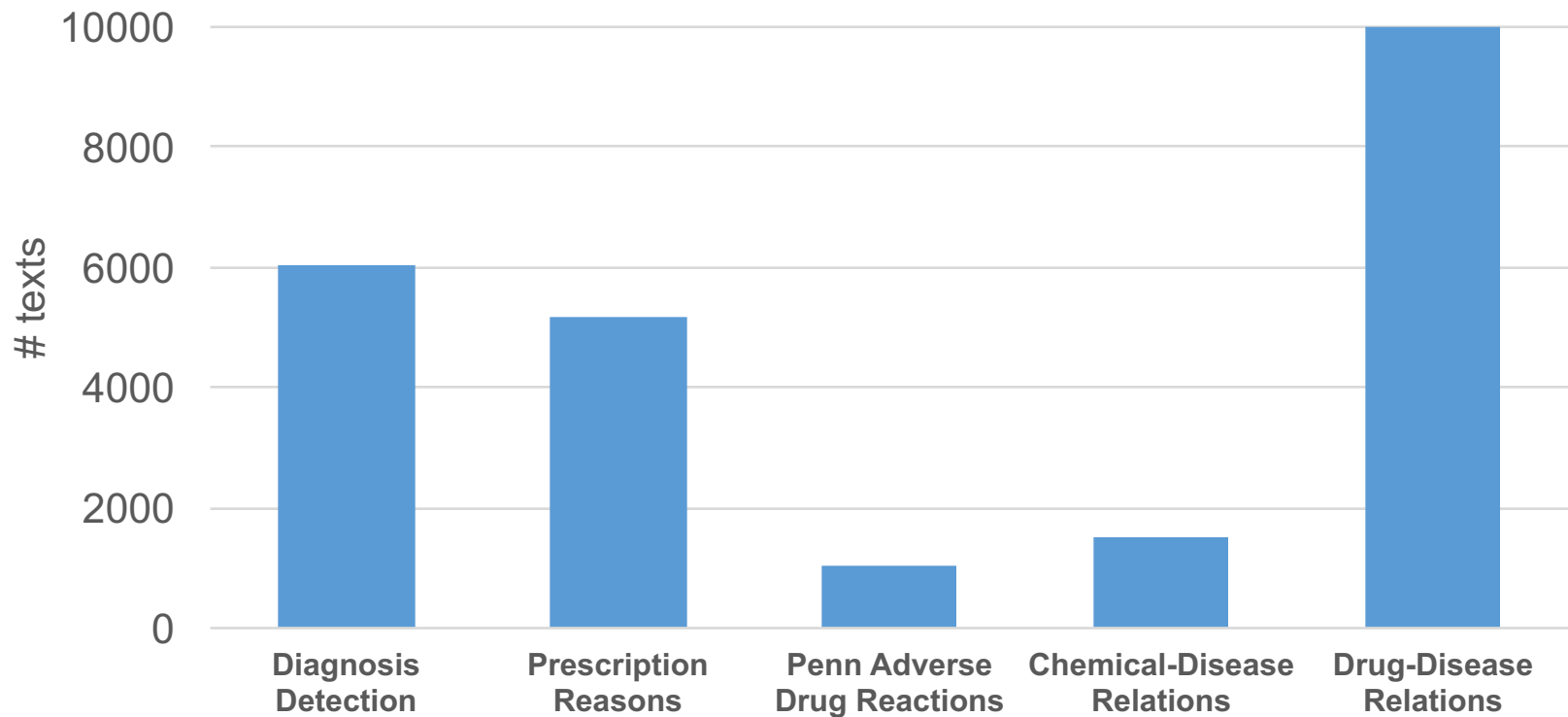
Commercial
Drug-Disease Relations



[Figures from: 1. Lamjed Ben Jabeur et al. Uprising microblogs: A Bayesian network retrieval model for tweet search. 2012. 2. <https://www.sjm.com.br/utilidades/pubmed-busca>. 3. <http://anakin.uta.cloud/uncategorized/the-need-for-drug-donations>. 4. <https://www.autismawareness.com.au/news-events/auupdate/is-there-an-over-diagnosis-of-autism/>]

Background: Healthcare Text Datasets

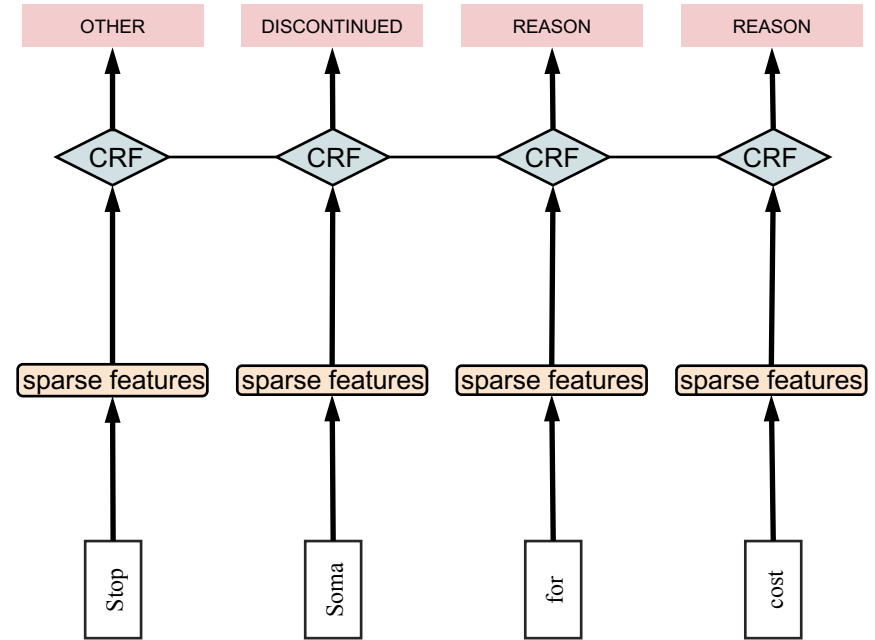
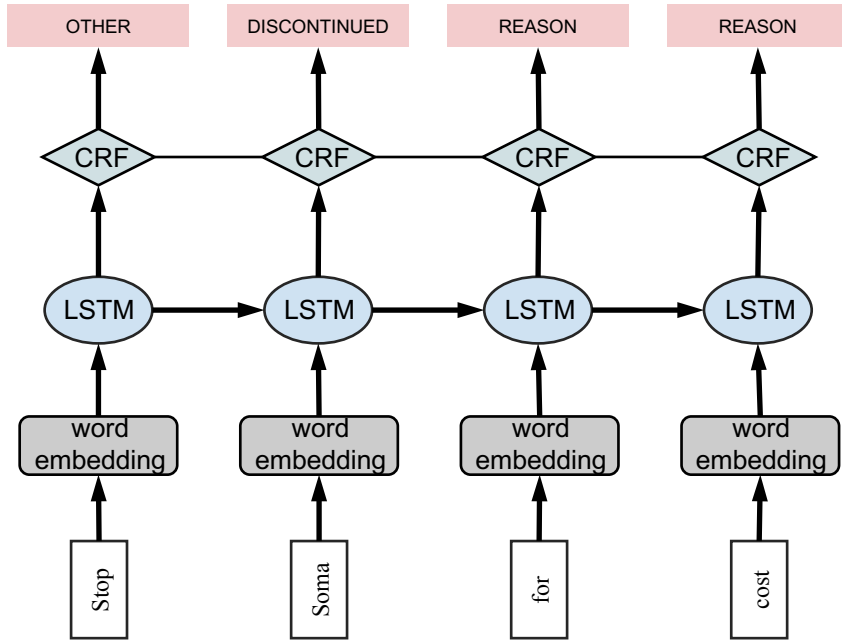
- Clinical text datasets are scarce and expensive
 - Privacy considerations
 - Domain specialists



Task: Clinical Text Annotation

Diagnosis Detection	<p>POSITIVE Asymptomatic bacteriuria , could be CONCERN neurogenic bladder disorder .</p>
Prescription Reasons	<p>PRESCRIBED REASON I will go ahead and place him on Clarinex for his seasonal allergic rhinitis .</p>
Penn Adverse Drug Reactions (ADR)	<p>ADR #TwoThingsThatDontMixWell venlafaxine and alcohol - you'll cry and ADR throw chairs at your mom's BBQ .</p>
Chemical–Disease Relations (CDR)	<p>DISEASE DRUG Ocular and auditory toxicity in hemodialyzed patients receiving desferrioxamine .</p>
Drug–Disease Relations	<p>TREATS Indicated for the management of active rheumatoid arthritis and should not be CONTRA used for rheumatoid arthritis in pregnant women .</p>

Previous Models



○ LSTM-CRF

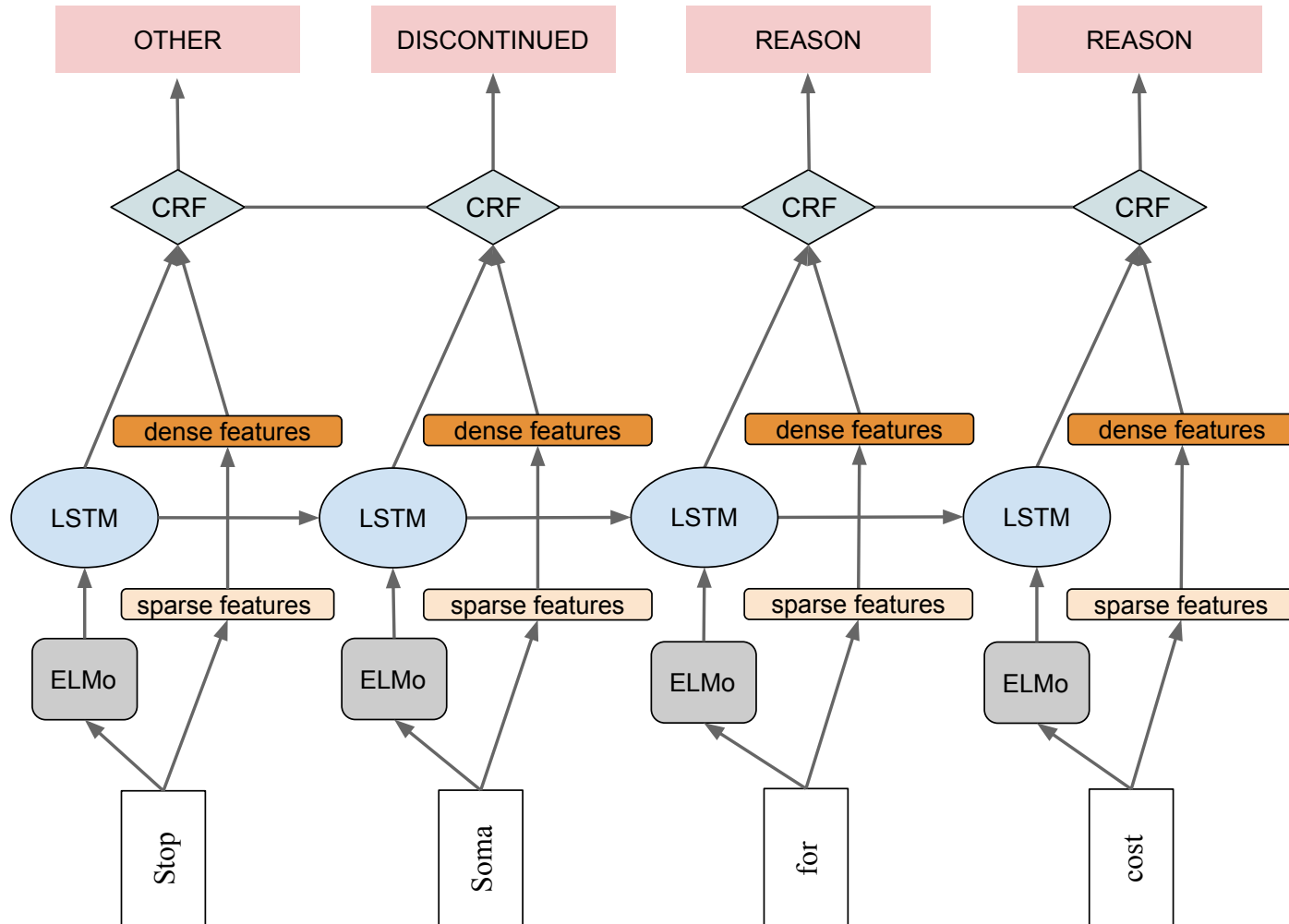
- General text
- Distributed word embeddings

○ HB-CRF

- Clinical text
- Sparse hand-built features

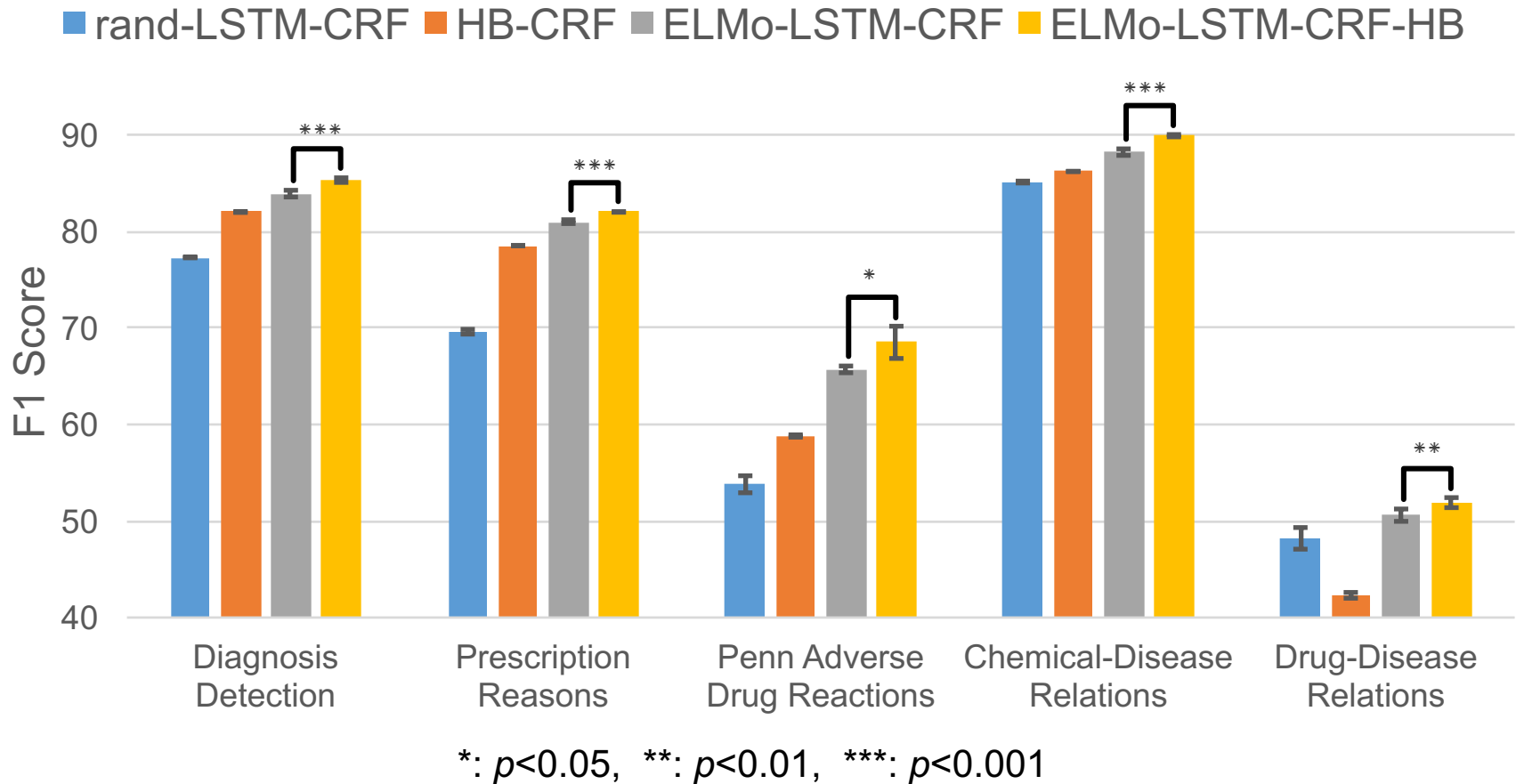
Model: ELMo-LSTM-CRF-HB

- Dense **ELMo** word embeddings + Sparse **hand-built** features



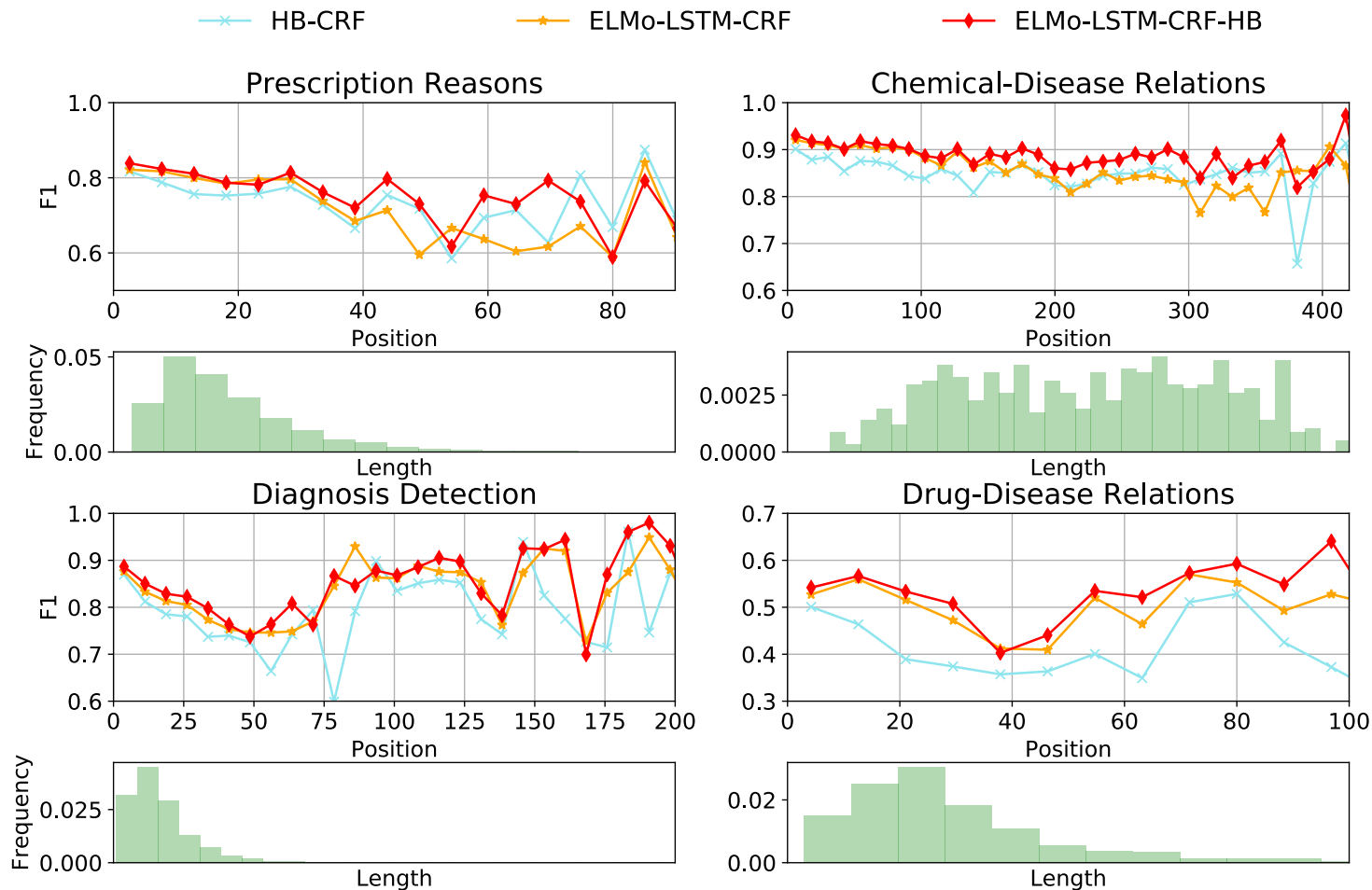
Performance: Per-token Macro-F1 Scores

- Hyperparameters tuned through cross-validation
- Each experiment repeated for five times

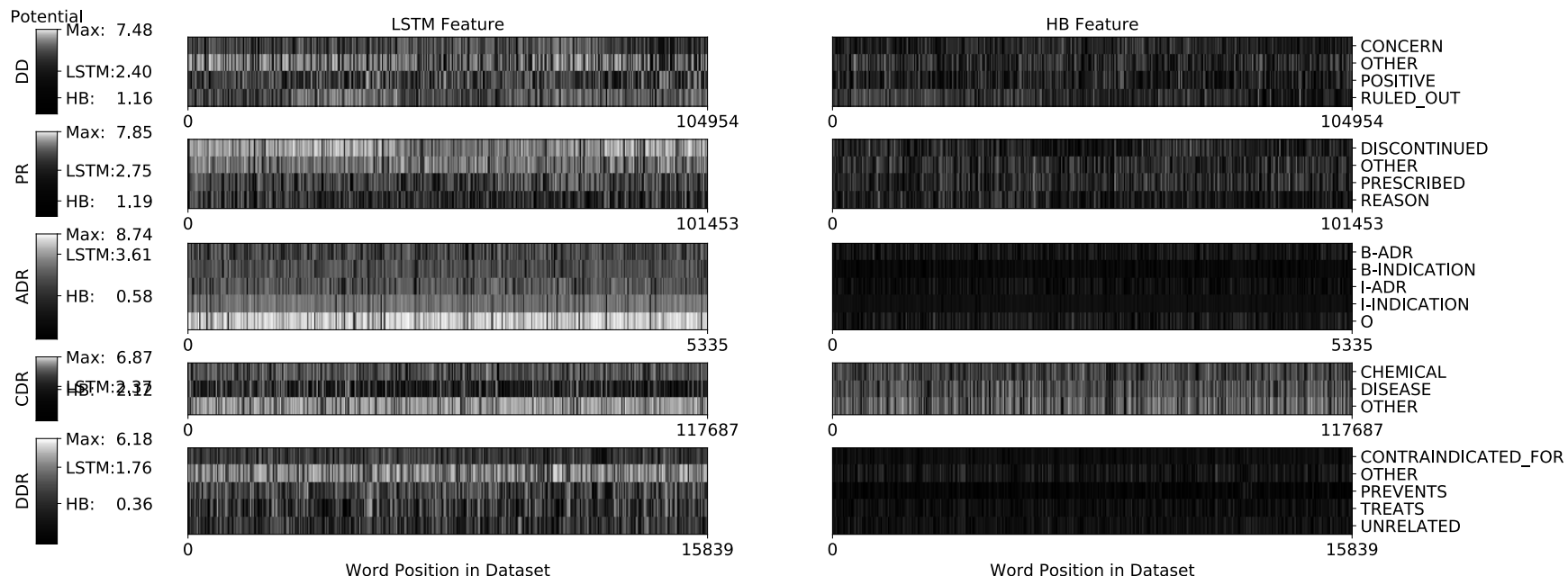


The Role of Text Length

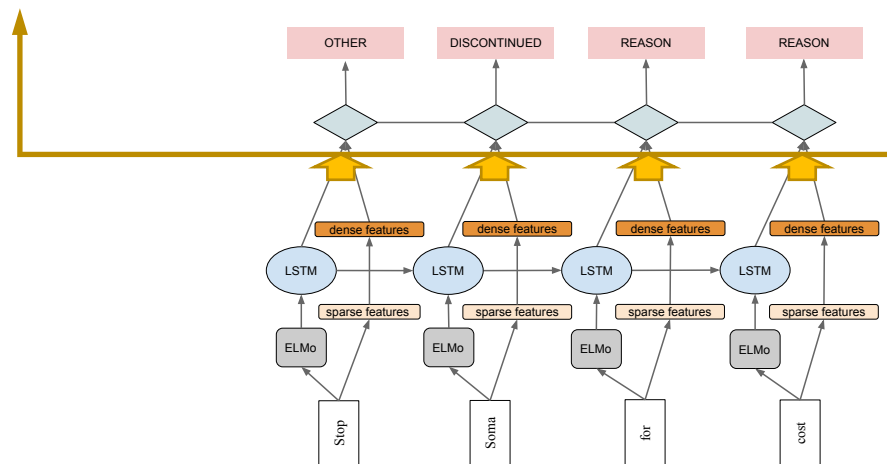
- LSTM: handles short texts well
- HB-CRF: robust on long texts



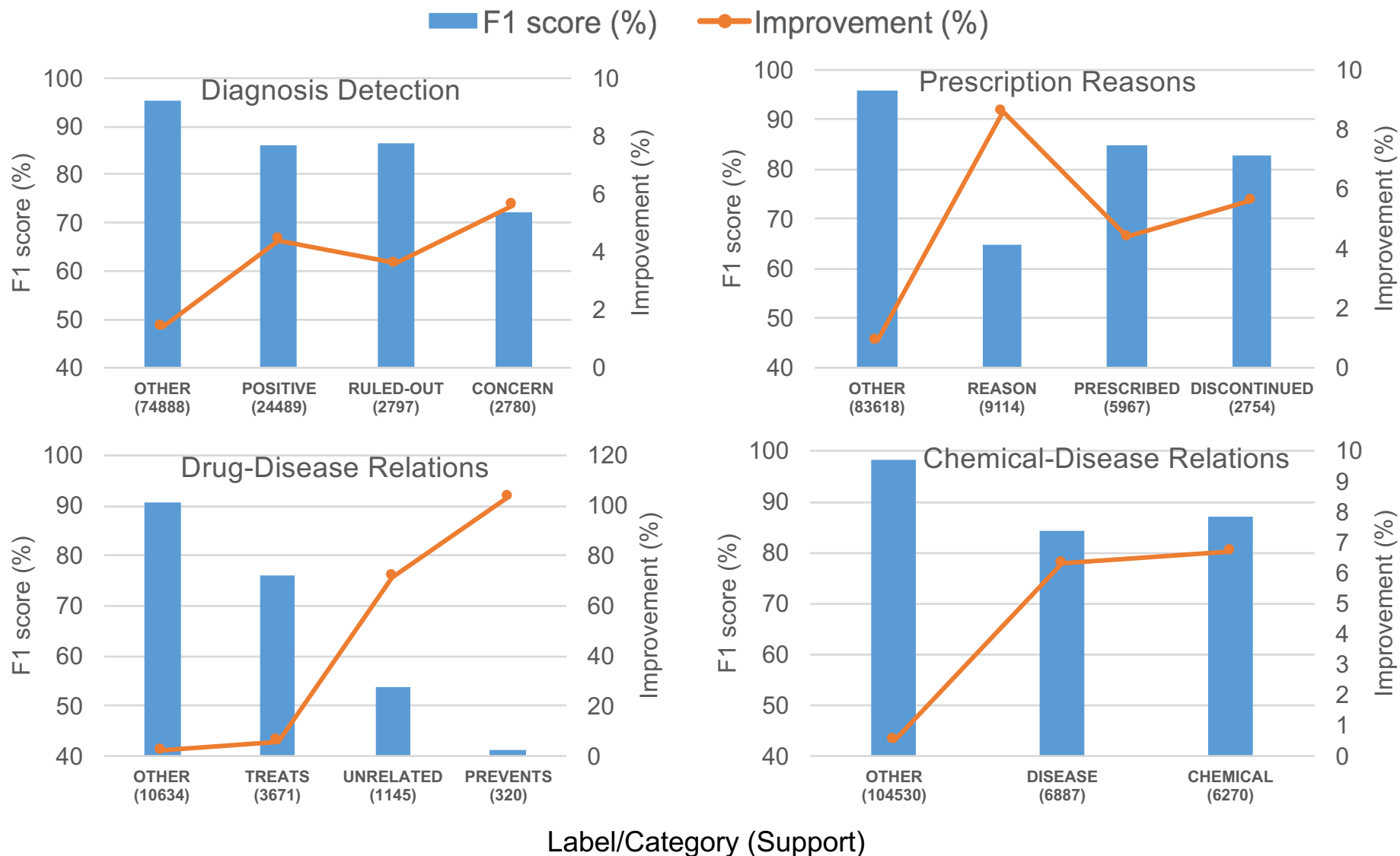
CRF Potential Scores



- LSTM features always more important
- HB features make substantial contribution



Major Improvements in Minor Categories



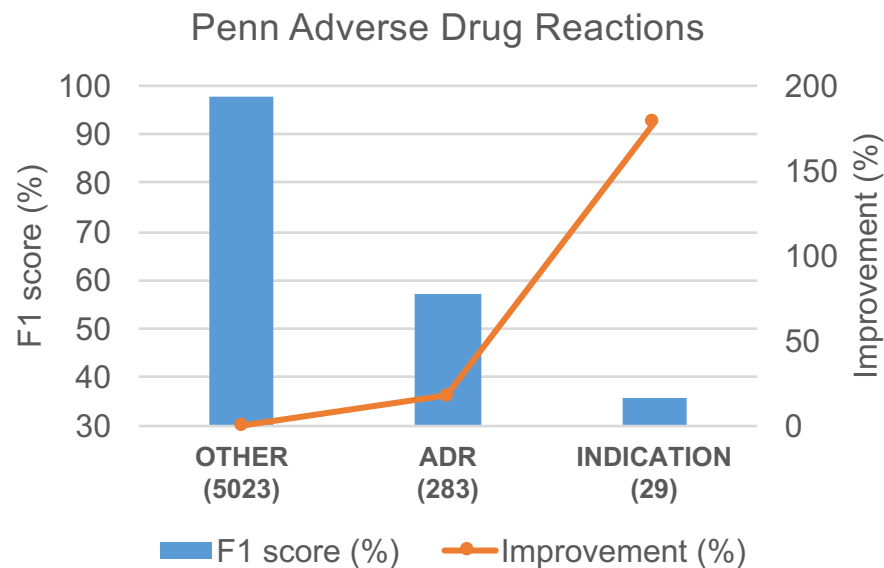
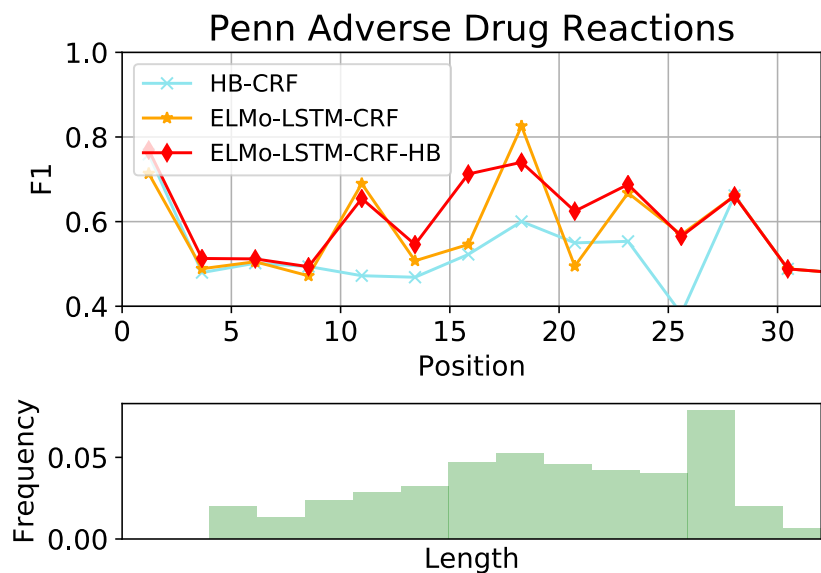
Conclusion

- A unified feature representation for clinical text sequence labeling
 - Sparse, ontology-driven features
 - Dense LSTM features
- Best performance on five distinct healthcare datasets
 - Takes advantages of both feature types
 - Makes maximal use of small, expensive, domain-specific healthcare texts
- A new labeled clinical dataset
 - Identifies the treatment relations between drugs and diseases
- Extensive analysis to identify what information our model makes use of, and why its performance is consistently improved

Acknowledgement

- Roam Analytics
 - Christopher Potts
 - Bruno Godefroy
 - Guillaume Genthial
 - Kevin Reschke
 - NLP Group

Penn Adverse Drug Reactions (ADR) Results

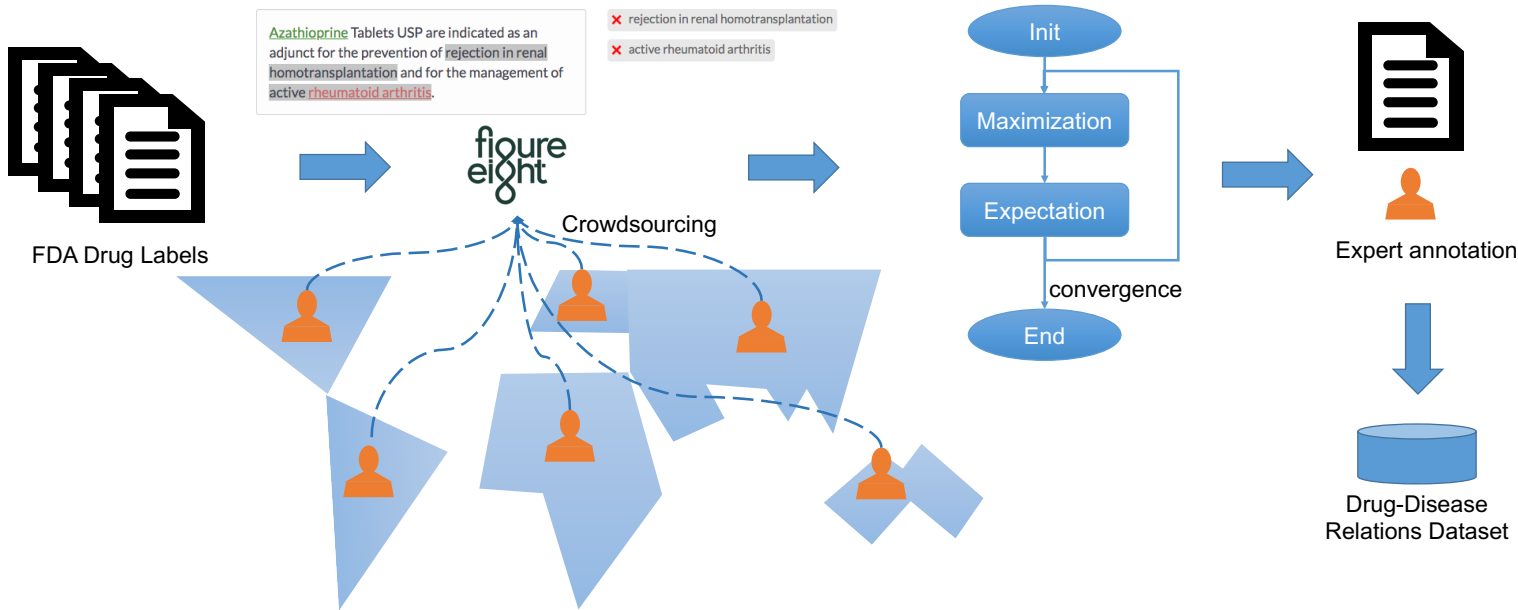


- The Role of Text Length
- Major Improvements in Minor Categories

Example of Hand-built Features

Sentence	Hand-built features of word <i>bacteria</i>
antiseptic handwash to decrease <i>bacteria</i> on the skin .	Adjacent words features: word-4:antiseptic, word-3:handwash, word-2:to, word-1:decrease, word:bacteria, word+1:on, word+2:the, word+3:skin, word+4:.. Adjacent POS tags features: tag-4:JJ, tag-3:NN, tag-2:TO, tag-1:VB, tag:NNS, tag+1:IN, tag+2:DT, tag+3:NN, tag+4:.. Semantic environment features: bias:1, is_upper:0, is_title:0, is_punctuation:0, in_left_context_of_negative_cues:0, in_right_context_of_negative_cues:0, in_left_context_of_prevents_cues:0, in_right_context_of_prevents_cues:0, in_left_context_of_treats_cues:0, in_right_context_of_treats_cues:0, in_left_context_of_treats_symptoms_cues:0, in_right_context_of_treats_symptoms_cues:0, in_left_context_of_contraindicated_cues:0, in_right_context_of_contraindicated_cues:0, in_left_context_of_affliction_adj_cues:0, in_right_context_of_affliction_adj_cues:0, in_left_context_of_indication_cues:0, in_right_context_of_indication_cues:0, in_left_context_of_details_cues:0, in_right_context_of_details_cues:0.

Procedure for Building Drug-Disease Relations Dataset



Statistics of Datasets

Statistics	Diagnosis Detection	Prescription Reasons	Penn Adverse Drug Reactions (ADR)	Chemical–Disease Relations (CDR)	Drug–Disease Relations
# texts	6042	5179	–	–	–
# training texts	–	–	749	1000	9494
# test texts	–	–	272	500	500
mean text length	17	19	19	227	30
max text length	374	258	40	623	542
# labels	4	4	5	3	5

Hyperparameters of Experiments

Models	Hyperparams	Diagnosis Detection	Prescription Reasons	Penn Adverse Drug Reactions (ADR)	Chemical–Disease Relations (CDR)	Drug–Disease Relations
rand-LSTM-CRF	η	1e-4	1e-4	1e-4	1e-4	1e-4
	epoch _{tune}	3	3	513	10	13
	epoch _{train}	34	40	3076	164	130
	\mathcal{R}_{c1}			{ 0, 3e-5, 1e-4, 3e-4, 1e-3 }		
	\mathcal{R}_{c2}			{ 0, 3e-4, 1e-3, 3e-3, 1e-2 }		
HB-CRF	η	1e-2	1e-2	3e-2	1e-2	1e-4
	epoch _{tune}	1	1	10	2	3
	epoch _{train}	3	4	82	10	35
	\mathcal{R}_{c1}			{ 0, 3e-6, 1e-5, 3e-5, 1e-4 }		
	\mathcal{R}_{c2}			{ 0, 3e-5, 1e-4, 3e-4, 1e-3 }		
ELMo-LSTM-CRF	η	1e-3	1e-3	1e-4	1e-3	5e-6
	epoch _{tune}	1	1	10	2	3
	epoch _{train}	3	4	82	10	35
	\mathcal{R}_{c1}			{ 0, 3e-5, 1e-4, 3e-4, 1e-3 }		
	\mathcal{R}_{c2}			{ 0, 3e-4, 1e-3, 3e-3, 1e-2 }		
ELMo-LSTM-CRF-HB	η	1e-3	1e-3	1e-4	1e-3	1e-5
	epoch _{tune}	1	1	10	2	3
	epoch _{train}	3	4	82	5	35
	\mathcal{R}_{c1}			{ 0, 3e-7, 1e-6, 3e-6, 1e-5 }		
	\mathcal{R}_{c2}			{ 0, 3e-6, 1e-5, 3e-5, 1e-4 }		