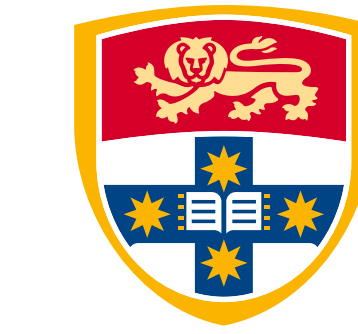


Recognizing Biomedical Names in Free Text

Xiang Dai

CSIRO Data61 and University of Sydney

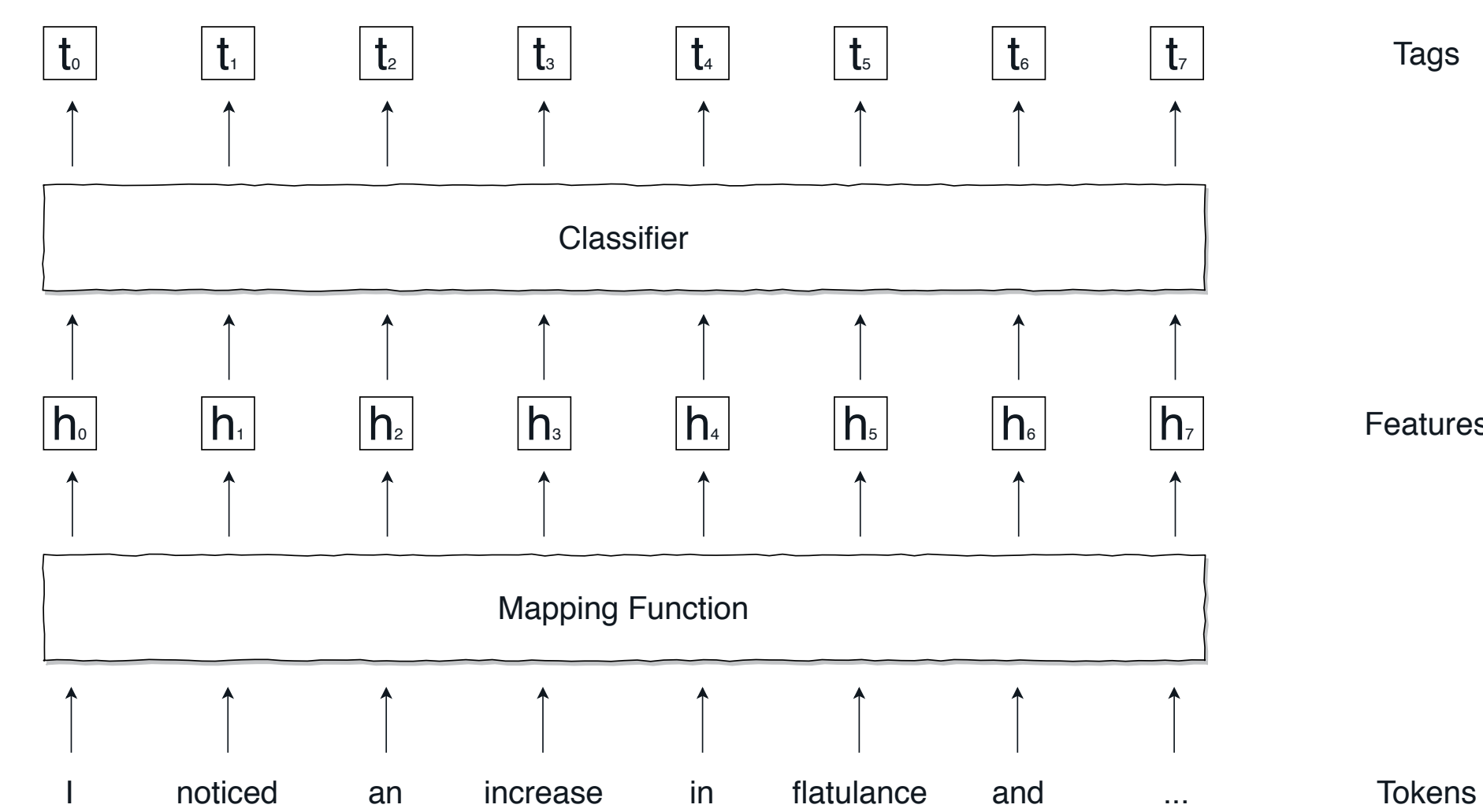


THE UNIVERSITY OF SYDNEY



Task and Prior Arts

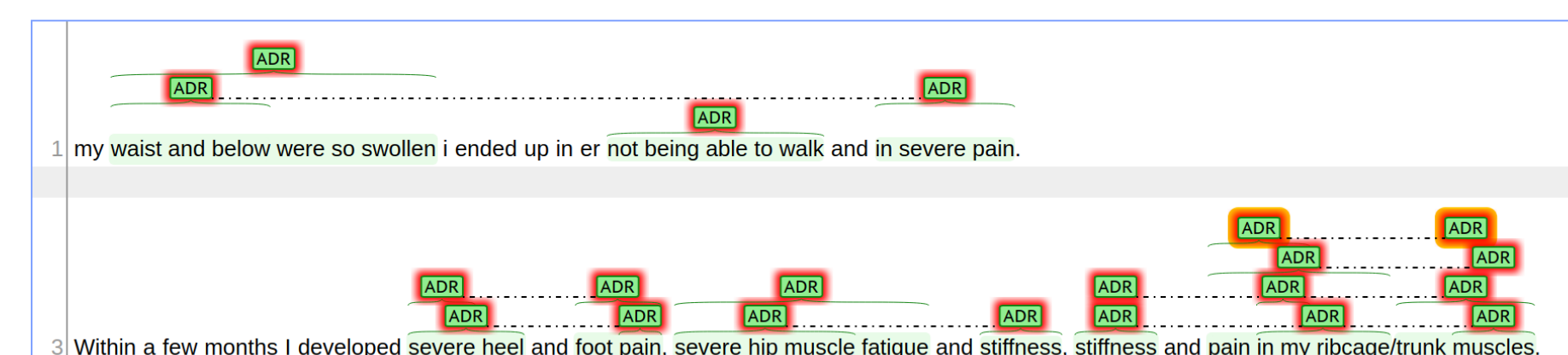
- The text, represented as a sequence of tokens, is taken as input, and entity mentions, each of which is represented as a set of token positions, are outputted. In addition, one entity category, such as disease, symptom, protein, DNA and so on, is assigned to each entity mention.
- The standard NER model – solving NER as a **sequence tagging** problem – assigns a tag to each token, indicating its role in the mention.



- Sequence taggers usually consist of two components: (1) mapping function converts the input sequence of tokens into a sequence of features vectors, each of which represents the corresponding token-in-the-context; (2) structural classifier predicts a sequence of tags given the input sequence of feature vectors.

Key Challenges

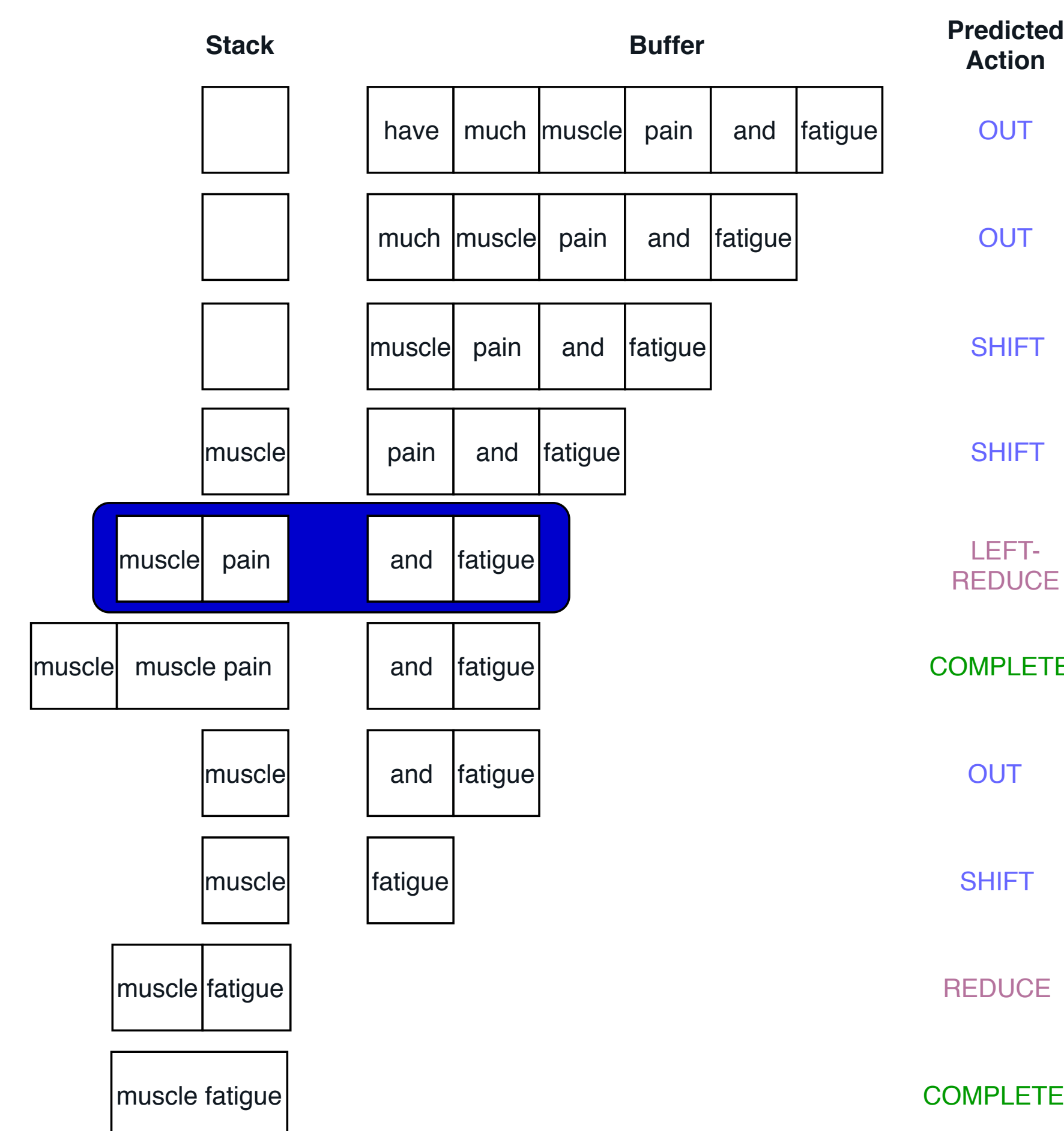
- Sequence tagger cannot recognize biomedical names that have **complex inner structures** – discontinuous and/or overlapping mentions.



- Training of deep neural models requires **large training set**, which is usually difficult to obtain in the biomedical domain.
- Discrepancy between **models pretrained on generic data** and biomedical data.

Transition-based model for discontinuous NER [1,2]

- The main motivation for recognizing discontinuous mentions is that they usually represent *compositional concepts* that differ from concepts represented by individual components.
- We propose a model based on the **shift-reduce parser** with six specialized actions and attention mechanism.
- The learning problem is then framed as: given the state of the parser, predict an action which is applied to change the state of the parser.



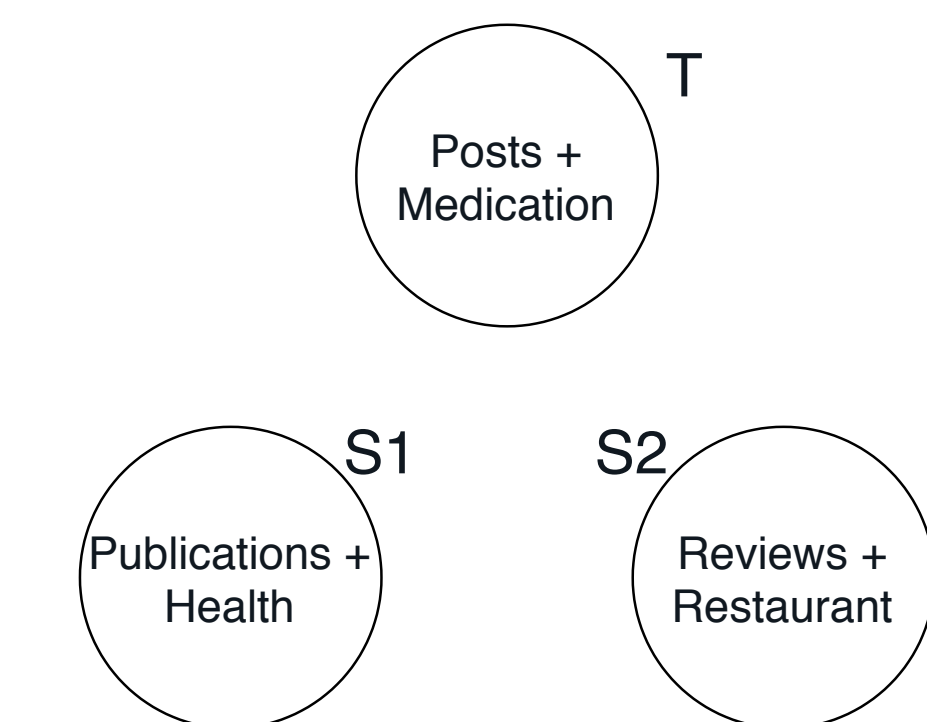
- Our model can effectively **recognize discontinuous mentions without sacrificing** the accuracy on continuous mentions.

References

- Dai, Xiang: Recognizing complex entity mentions: A review and future directions. In SRW@ACL; 2018: 37-44.
- Dai, Xiang, Karimi, Sarvnaz, Hachey, Ben, Paris, Cecile: An Effective Transition-based Model for Discontinuous NER. In ACL; 2020: 5860-5870.
- Dai, Xiang, Karimi, Sarvnaz, Hachey, Ben, Paris, Cecile: Using similarity measures to select pretraining data for NER. In NAACL; 2019: 1460-1470.
- Dai, Xiang, Karimi, Sarvnaz, Hachey, Ben, Paris, Cecile: Cost-effective Selection of Pretraining Data: A Case Study of Pretraining BERT on Social Media. In Findings of EMNLP; 2020: 1675-1681.
- Dai, Xiang, Adel, Heike: An Analysis of Simple Data Augmentation for Named Entity Recognition. In COLING; 2020: 3861-3867.

Select in-domain Pretrain Data [3,4]

- Studies on domain-specific pretrained models show that, when in-domain data is used for pretraining, target task performance can be improved. However, the selection of in-domain data usually resorts to **human intuition**.
- We find that human intuition varies across NLP practitioners, especially regarding intersecting domains.



- We employ **cost-effective** measures to quantify similarity between source pretraining and target task data, and demonstrate that these measures are good predictors of the usefulness of pretrained models.

Data augmentation for NER [5]

- We design several **easy to adapt** data augmentation for NER, and show that simple data augmentation can improve performance even over strong baselines.

	Instance							
None	She	did	not	complain	of	headache	or	
	O	O	O	O	O	B-problem	O	
	any	other	neurological	symptoms	.			
	B-problem	I-problem	I-problem	I-problem	O			
LWTR	L	One	not	complain	of	headache	he	
	O	O	O	O	O	B-problem	O	
	any	interatrial	neurological	current	.			
	B-problem	I-problem	I-problem	I-problem	O			
SR	She	did	non	complain	of	headache	or	
	O	O	O	O	O	B-problem	O	
	whatsoever	former	neurologic	symptom	.			
	B-problem	I-problem	I-problem	I-problem	O			
MR	She	did	not	complain	of	neuropathic	pain	
	O	O	O	O	O	B-problem	I-problem	
	syndrome	or	acute	pulmonary	disease	.		
	I-problem	O	B-problem	I-problem	I-problem	O		
SIS	not	complain	She	did	of	headache	or	
	O	O	O	O	O	B-problem	O	
	neurological	any	symptoms	other	.			
	B-problem	I-problem	I-problem	I-problem	O			