

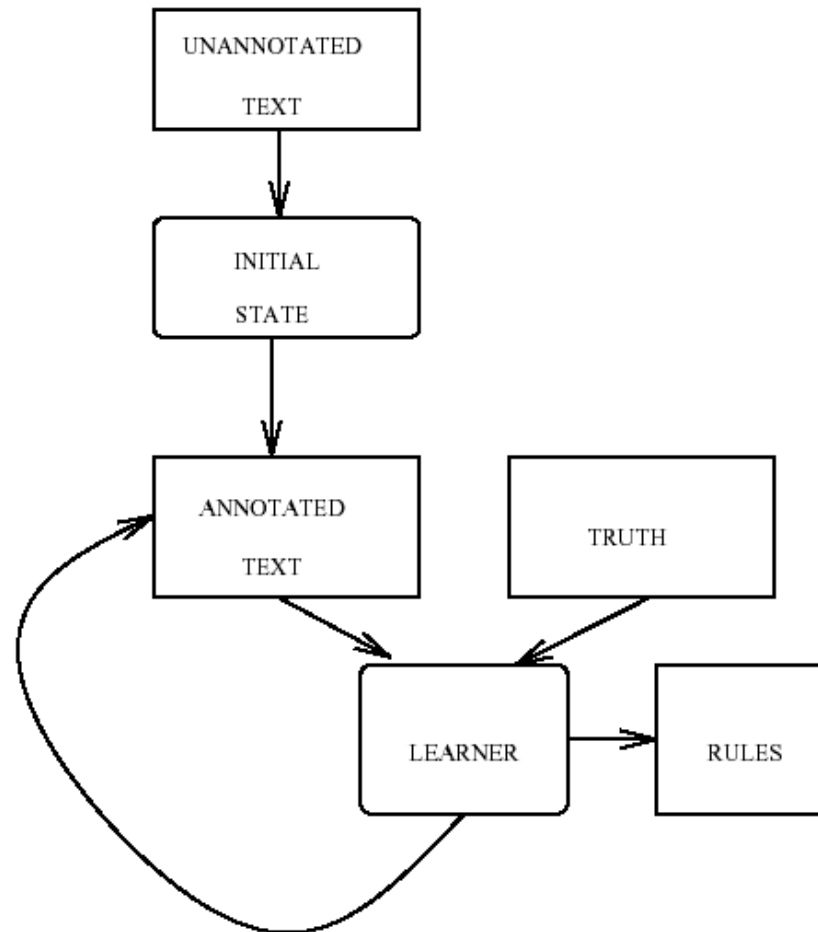
# CMPT-825

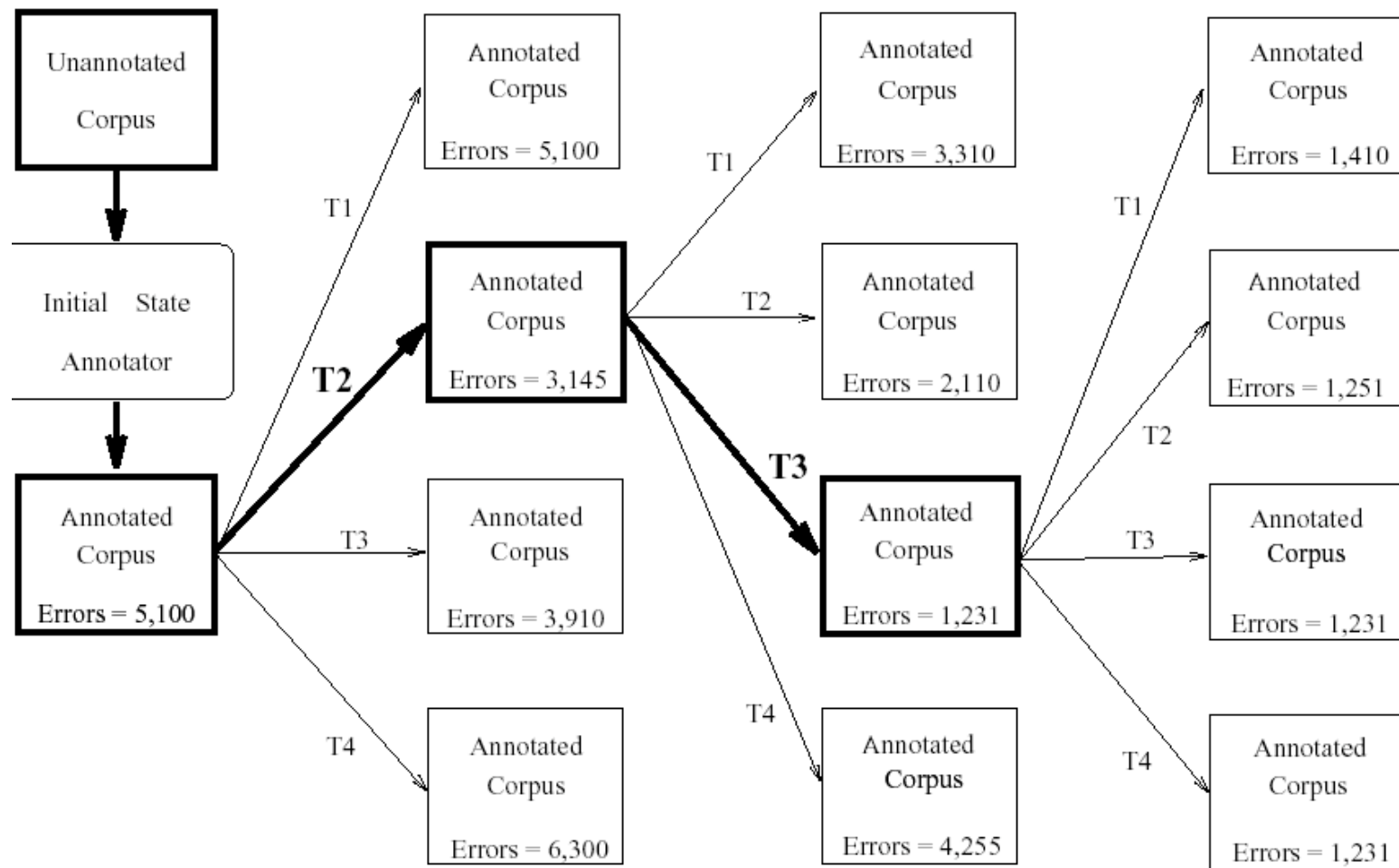
## Natural Language Processing

Anoop Sarkar

`http://www.cs.sfu.ca/~anoop`

## Transformation-Based Error-Driven Learning





1. apply initial state annotator to corpus
2. **while** transformations can still be found **do**
3.   **for** `from_tag = tag1` **to** `tagn`
4.     **for** `to_tag = tag1` **to** `tagn`
5.       **for** `corpus_position = 1` **to** `corpus_size`
6.       **if** (`correct_tag(corpus_position) == to_tag`  
           && `current_tag(corpus_position) == from_tag`)
7.           `num_good_transformations(tag(corpus_position - 1))++`
8.       **else if** (`correct_tag(corpus_position) == from_tag`  
           && `current_tag(corpus_position) == from_tag`)
9.           `num_bad_transformations(tag(corpus_position - 1))++`
10.    find  $max_T$  (`num_good_transformations(T) - num_bad_transformations(T)`)
11.    **if** this is the best scoring rule found yet then store as best rule:  
       Change tag from `from_tag` to `to_tag` if previous tag is T
12. apply best rule to training corpus
13. append best rule to ordered list of transformations

Change Tag			
#	From	To	Condition
1	NN	VB	Previous tag is <i>TO</i>
2	VBP	VB	One of the previous three tags is <i>MD</i>
3	NN	VB	One of the previous two tags is <i>MD</i>
4	VB	NN	One of the previous two tags is <i>DT</i>
5	VBD	VBN	One of the previous three tags is <i>VBZ</i>
6	VBN	VBD	Previous tag is <i>PRP</i>
7	VBN	VBD	Previous tag is <i>NNP</i>
8	VBD	VBN	Previous tag is <i>VBD</i>
9	VBP	VB	Previous tag is <i>TO</i>
10	POS	VBZ	Previous tag is <i>PRP</i>
11	VB	VBP	Previous tag is <i>NNS</i>
12	VBD	VBN	One of previous three tags is <i>VBP</i>
13	IN	WDT	One of next two tags is <i>VB</i>
14	VBD	VBN	One of previous two tags is <i>VB</i>
15	VB	VBP	Previous tag is <i>PRP</i>
16	IN	WDT	Next tag is <i>VBZ</i>
17	IN	DT	Next tag is <i>NN</i>
18	JJ	NNP	Next tag is <i>NNP</i>
19	IN	WDT	Next tag is <i>VBD</i>
20	JJR	RBR	Next tag is <i>JJ</i>

Change tag **a** to tag **b** when:

1. The preceding (following) word is  $w$ .
2. The word two before (after) is  $w$ .
3. One of the two preceding (following) words is  $w$ .
4. The current word is  $w$  and the preceding (following) word is  $x$ .
5. The current word is  $w$  and the preceding (following) word is tagged  $z$ .
6. The current word is  $w$ .
7. The preceding (following) word is  $w$  and the preceding (following) tag is  $t$ .
8. The current word is  $w$ , the preceding (following) word is  $w_2$  and the preceding (following) tag is  $t$ .

Method	Training Corpus Size (Words)	# of Rules or Context. Probs.	Acc. (%)
Stochastic	64 K	6,170	96.3
Stochastic	1 Million	10,000	96.7
Rule-Based With Lex. Rules	64 K	215	96.7
Rule-Based With Lex. Rules	600 K	447	97.2
Rule-Based w/o Lex. Rules	600 K	378	97.0