

Supertagging in PET

Rebecca Dridan

Universität des Saarlandes

DELPH-IN Summit
Barcelona, July 2009

Outline

What is supertagging?

Modifying the grammar

How does it work?

Configuration options

What are the supertags?

What is supertagging?

Supertagging is restricting the parser search space by limiting the number of lexical items, in order to increase parser efficiency.

CCG lexicon entries

```
price ((S[b]\NP)/PP)/NP
price (S[b]\NP)/NP
price (S[dcl]\NP)/NP
price N
price N/N
price NP\NP
```

HPSG lexicon entry

```
price_n2 := n_pp_c-obl_le &
```

```
[STEM < "price" >,
 [SYNSEM [LKEYS [-COMPKEY _of_p_sel_rel
                 KEYREL.PRED "_price_n_of_rel"],
          PHON [ONSET con]]]]
```

In CCG, all information is in the supertag (category) and the word form → the supertagger produces the lexical items.

In HPSG, this is only true if the supertag is the lexical entry (**price_n2**) → the supertagger is a filter for the lexicon.

Modifying the grammar

Supertagging is facilitated by the new `TOKENS` feature of `word_or_lexrule`, added for chart mapping.

Step 1 Add `+STAG` feature to the **token** type:

token

<table style="width: 100%; border-collapse: collapse;"> <tr><td style="padding: 2px 10px;"><code>+FORM</code></td><td style="padding: 2px 10px;"><code>string</code></td></tr> <tr><td style="padding: 2px 10px;"><code>+CLASS</code></td><td style="padding: 2px 10px;"><code>token_class</code></td></tr> <tr><td style="padding: 2px 10px;"><code>+TRAIT</code></td><td style="padding: 2px 10px;"><code>token_trait</code></td></tr> <tr><td style="padding: 2px 10px;"><code>+PRED</code></td><td style="padding: 2px 10px;"><code>predsort</code></td></tr> <tr><td style="padding: 2px 10px;"><code>+CARG</code></td><td style="padding: 2px 10px;"><code>string</code></td></tr> <tr><td style="padding: 2px 10px;"><code>+ID</code></td><td style="padding: 2px 10px;"><code>*diff-list*</code></td></tr> <tr><td style="padding: 2px 10px;"><code>+FROM</code></td><td style="padding: 2px 10px;"><code>string</code></td></tr> <tr><td style="padding: 2px 10px;"><code>+TO</code></td><td style="padding: 2px 10px;"><code>string</code></td></tr> <tr><td style="padding: 2px 10px;"><code>+TNT</code></td><td style="padding: 2px 10px;"> <table style="border-collapse: collapse;"> <tr><td style="padding: 2px 10px;"><code>+TAGS</code></td><td style="padding: 2px 10px;"><code>*list*</code></td></tr> <tr><td style="padding: 2px 10px;"><code>+PRBS</code></td><td style="padding: 2px 10px;"><code>*list*</code></td></tr> </table> </td></tr> </table>	<code>+FORM</code>	<code>string</code>	<code>+CLASS</code>	<code>token_class</code>	<code>+TRAIT</code>	<code>token_trait</code>	<code>+PRED</code>	<code>predsort</code>	<code>+CARG</code>	<code>string</code>	<code>+ID</code>	<code>*diff-list*</code>	<code>+FROM</code>	<code>string</code>	<code>+TO</code>	<code>string</code>	<code>+TNT</code>	<table style="border-collapse: collapse;"> <tr><td style="padding: 2px 10px;"><code>+TAGS</code></td><td style="padding: 2px 10px;"><code>*list*</code></td></tr> <tr><td style="padding: 2px 10px;"><code>+PRBS</code></td><td style="padding: 2px 10px;"><code>*list*</code></td></tr> </table>	<code>+TAGS</code>	<code>*list*</code>	<code>+PRBS</code>	<code>*list*</code>	+	<table style="border-collapse: collapse;"> <tr> <td style="padding: 5px;"><code>+STAG</code></td> <td style="padding: 5px;"> <table style="border-collapse: collapse;"> <tr> <td style="padding: 5px;"><code>+STAGS</code></td> <td style="padding: 5px;"><code>*list*</code></td> </tr> <tr> <td style="padding: 5px;"><code>+SPRBS</code></td> <td style="padding: 5px;"><code>*list*</code></td> </tr> </table> </td> </tr> <tr> <td style="padding: 5px;"><code>stag</code></td> <td style="padding: 5px;"></td> </tr> </table>	<code>+STAG</code>	<table style="border-collapse: collapse;"> <tr> <td style="padding: 5px;"><code>+STAGS</code></td> <td style="padding: 5px;"><code>*list*</code></td> </tr> <tr> <td style="padding: 5px;"><code>+SPRBS</code></td> <td style="padding: 5px;"><code>*list*</code></td> </tr> </table>	<code>+STAGS</code>	<code>*list*</code>	<code>+SPRBS</code>	<code>*list*</code>	<code>stag</code>	
<code>+FORM</code>	<code>string</code>																															
<code>+CLASS</code>	<code>token_class</code>																															
<code>+TRAIT</code>	<code>token_trait</code>																															
<code>+PRED</code>	<code>predsort</code>																															
<code>+CARG</code>	<code>string</code>																															
<code>+ID</code>	<code>*diff-list*</code>																															
<code>+FROM</code>	<code>string</code>																															
<code>+TO</code>	<code>string</code>																															
<code>+TNT</code>	<table style="border-collapse: collapse;"> <tr><td style="padding: 2px 10px;"><code>+TAGS</code></td><td style="padding: 2px 10px;"><code>*list*</code></td></tr> <tr><td style="padding: 2px 10px;"><code>+PRBS</code></td><td style="padding: 2px 10px;"><code>*list*</code></td></tr> </table>	<code>+TAGS</code>	<code>*list*</code>	<code>+PRBS</code>	<code>*list*</code>																											
<code>+TAGS</code>	<code>*list*</code>																															
<code>+PRBS</code>	<code>*list*</code>																															
<code>+STAG</code>	<table style="border-collapse: collapse;"> <tr> <td style="padding: 5px;"><code>+STAGS</code></td> <td style="padding: 5px;"><code>*list*</code></td> </tr> <tr> <td style="padding: 5px;"><code>+SPRBS</code></td> <td style="padding: 5px;"><code>*list*</code></td> </tr> </table>	<code>+STAGS</code>	<code>*list*</code>	<code>+SPRBS</code>	<code>*list*</code>																											
<code>+STAGS</code>	<code>*list*</code>																															
<code>+SPRBS</code>	<code>*list*</code>																															
<code>stag</code>																																

Modifying the grammar

Step 2 For each lexical type that may be filtered, specify the +STAG feature.

for example:

v_np_le :+

$$\left[\text{TOKENS.LIST} \quad \left\langle \left[+\text{STAG.}+\text{STAGS} \quad \left\langle \text{"v_np_le", \dots} \right\rangle \right] \right\rangle \right]$$

This is the first point of supertagger control: if a particular lexical type should never be filtered by the supertagger, leave the +STAG feature unspecified.

Note: the value is of type string and may be set to any string value.

Modifying the grammar

Step 3 Modify token mapping rules

- Make sure +STAG information is copied through where appropriate.
- Expand single token with multiple supertags into multiple tokens each with one supertag?
- Delete +STAG according the tag probabilities?
- Deal with TnT POS tag appropriately.
 - copy?
 - delete?
 - modify?

This is the second point of control that dictates how the supertags will be used in parsing.

How does it work?

lexical entries

<table style="border-collapse: collapse; width: 100%;"> <tr> <td style="padding: 5px;">STEM</td> <td style="padding: 5px;">⟨ "ration" ⟩</td> </tr> <tr> <td style="padding: 5px;">SYNSEM</td> <td style="padding: 5px;"> <table style="border-collapse: collapse; width: 100%;"> <tr> <td style="padding: 5px;">KEYREL.PRED</td> <td style="padding: 5px;">"_ration_v_1_rel"</td> </tr> <tr> <td style="padding: 5px;">PHON.ONSET</td> <td style="padding: 5px;">con</td> </tr> </table> </td> </tr> <tr> <td style="padding: 5px;">TOKENS</td> <td style="padding: 5px;">⟨ 1 [+STAG.+STAGS ⟨ "vp_np_le" ⟩] ⟩</td> </tr> </table>	STEM	⟨ "ration" ⟩	SYNSEM	<table style="border-collapse: collapse; width: 100%;"> <tr> <td style="padding: 5px;">KEYREL.PRED</td> <td style="padding: 5px;">"_ration_v_1_rel"</td> </tr> <tr> <td style="padding: 5px;">PHON.ONSET</td> <td style="padding: 5px;">con</td> </tr> </table>	KEYREL.PRED	"_ration_v_1_rel"	PHON.ONSET	con	TOKENS	⟨ 1 [+STAG.+STAGS ⟨ "vp_np_le" ⟩] ⟩	<table style="border-collapse: collapse; width: 100%;"> <tr> <td style="padding: 5px;">STEM</td> <td style="padding: 5px;">⟨ "ration" ⟩</td> </tr> <tr> <td style="padding: 5px;">SYNSEM</td> <td style="padding: 5px;"> <table style="border-collapse: collapse; width: 100%;"> <tr> <td style="padding: 5px;">KEYREL.PRED</td> <td style="padding: 5px;">"_ration_n_1_rel"</td> </tr> <tr> <td style="padding: 5px;">PHON.ONSET</td> <td style="padding: 5px;">con</td> </tr> </table> </td> </tr> <tr> <td style="padding: 5px;">TOKENS</td> <td style="padding: 5px;">⟨ 1 [+STAG.+STAGS ⟨ "n_-_c_le" ⟩] ⟩</td> </tr> </table>	STEM	⟨ "ration" ⟩	SYNSEM	<table style="border-collapse: collapse; width: 100%;"> <tr> <td style="padding: 5px;">KEYREL.PRED</td> <td style="padding: 5px;">"_ration_n_1_rel"</td> </tr> <tr> <td style="padding: 5px;">PHON.ONSET</td> <td style="padding: 5px;">con</td> </tr> </table>	KEYREL.PRED	"_ration_n_1_rel"	PHON.ONSET	con	TOKENS	⟨ 1 [+STAG.+STAGS ⟨ "n_-_c_le" ⟩] ⟩
STEM	⟨ "ration" ⟩																				
SYNSEM	<table style="border-collapse: collapse; width: 100%;"> <tr> <td style="padding: 5px;">KEYREL.PRED</td> <td style="padding: 5px;">"_ration_v_1_rel"</td> </tr> <tr> <td style="padding: 5px;">PHON.ONSET</td> <td style="padding: 5px;">con</td> </tr> </table>	KEYREL.PRED	"_ration_v_1_rel"	PHON.ONSET	con																
KEYREL.PRED	"_ration_v_1_rel"																				
PHON.ONSET	con																				
TOKENS	⟨ 1 [+STAG.+STAGS ⟨ "vp_np_le" ⟩] ⟩																				
STEM	⟨ "ration" ⟩																				
SYNSEM	<table style="border-collapse: collapse; width: 100%;"> <tr> <td style="padding: 5px;">KEYREL.PRED</td> <td style="padding: 5px;">"_ration_n_1_rel"</td> </tr> <tr> <td style="padding: 5px;">PHON.ONSET</td> <td style="padding: 5px;">con</td> </tr> </table>	KEYREL.PRED	"_ration_n_1_rel"	PHON.ONSET	con																
KEYREL.PRED	"_ration_n_1_rel"																				
PHON.ONSET	con																				
TOKENS	⟨ 1 [+STAG.+STAGS ⟨ "n_-_c_le" ⟩] ⟩																				

are unified with

1	<table style="border-collapse: collapse; width: 100%;"> <tr> <td style="padding: 5px;">+FORM</td> <td style="padding: 5px;">"rations"</td> </tr> <tr> <td style="padding: 5px;">+FROM</td> <td style="padding: 5px;">"11"</td> </tr> <tr> <td style="padding: 5px;">+TO</td> <td style="padding: 5px;">"18"</td> </tr> <tr> <td style="padding: 5px;">+TNT.+TAGS</td> <td style="padding: 5px;">⟨ "NNS" ⟩</td> </tr> <tr> <td style="padding: 5px;">+TNT.+PRBS</td> <td style="padding: 5px;">⟨ "0.97305" ⟩</td> </tr> <tr> <td style="padding: 5px;">+STAG.+STAGS</td> <td style="padding: 5px;">⟨ "n_-_c_le" ⟩</td> </tr> <tr> <td style="padding: 5px;">+STAG.+SPRBS</td> <td style="padding: 5px;">⟨ "0.60377" ⟩</td> </tr> </table>	+FORM	"rations"	+FROM	"11"	+TO	"18"	+TNT.+TAGS	⟨ "NNS" ⟩	+TNT.+PRBS	⟨ "0.97305" ⟩	+STAG.+STAGS	⟨ "n_-_c_le" ⟩	+STAG.+SPRBS	⟨ "0.60377" ⟩
+FORM	"rations"														
+FROM	"11"														
+TO	"18"														
+TNT.+TAGS	⟨ "NNS" ⟩														
+TNT.+PRBS	⟨ "0.97305" ⟩														
+STAG.+STAGS	⟨ "n_-_c_le" ⟩														
+STAG.+SPRBS	⟨ "0.60377" ⟩														

the input token

How does it work?



INVALID

STEM	⟨ "ration" ⟩														
SYNSEM	<table style="border-collapse: collapse;"> <tr> <td style="border-right: 1px solid black; padding: 5px;">KEYREL.PRED</td> <td style="padding: 5px;">"_ration_n_1_rel"</td> </tr> <tr> <td style="border-right: 1px solid black; padding: 5px;">PHON.ONSET</td> <td style="padding: 5px;">con</td> </tr> </table>	KEYREL.PRED	"_ration_n_1_rel"	PHON.ONSET	con										
KEYREL.PRED	"_ration_n_1_rel"														
PHON.ONSET	con														
TOKENS	<table style="border-collapse: collapse;"> <tr> <td style="border-right: 1px solid black; padding: 5px;">+FORM</td> <td style="padding: 5px;">"rations"</td> </tr> <tr> <td style="border-right: 1px solid black; padding: 5px;">+FROM</td> <td style="padding: 5px;">"11"</td> </tr> <tr> <td style="border-right: 1px solid black; padding: 5px;">+TO</td> <td style="padding: 5px;">"18"</td> </tr> <tr> <td style="border-right: 1px solid black; padding: 5px;">+TNT.+TAGS</td> <td style="padding: 5px;">⟨ "NNS" ⟩</td> </tr> <tr> <td style="border-right: 1px solid black; padding: 5px;">+TNT.+PRBS</td> <td style="padding: 5px;">⟨ "0.97305" ⟩</td> </tr> <tr> <td style="border-right: 1px solid black; padding: 5px;">+STAG.+STAGS</td> <td style="padding: 5px;">⟨ "n_-_c_1e" ⟩</td> </tr> <tr> <td style="border-right: 1px solid black; padding: 5px;">+STAG.+SPRBS</td> <td style="padding: 5px;">⟨ "0.60377" ⟩</td> </tr> </table>	+FORM	"rations"	+FROM	"11"	+TO	"18"	+TNT.+TAGS	⟨ "NNS" ⟩	+TNT.+PRBS	⟨ "0.97305" ⟩	+STAG.+STAGS	⟨ "n_-_c_1e" ⟩	+STAG.+SPRBS	⟨ "0.60377" ⟩
+FORM	"rations"														
+FROM	"11"														
+TO	"18"														
+TNT.+TAGS	⟨ "NNS" ⟩														
+TNT.+PRBS	⟨ "0.97305" ⟩														
+STAG.+STAGS	⟨ "n_-_c_1e" ⟩														
+STAG.+SPRBS	⟨ "0.60377" ⟩														

Configuration options

Handling lexical gaps

When the supertagger filters out every lexical entry, do we

- **fail to parse?**
Supertagging is often a trade-off between robustness and efficiency. An option when the supertagger is very good, or efficiency has much higher priority than robustness.
- **create a generic item from the supertag?**
Similar effect to CCG supertagging.
- **use the TnT POS tag to generate a generic lexical item?**
This is the standard method of handling any lexical gap. Should compatibility between supertag and POS tag be enforced?

Configuration options

Assigning supertags

Which input tokens should be restricted, and how?

- **single tag all tokens**
Simple, but brutal. Best sentence tagging accuracy approximately 50% → low coverage, high efficiency.
- **multitag all tokens**
Optimal β value must be determined. Sentence tagging accuracy of 91% possible assigning an average of 2.3 tags per token.
- **tag only some tokens, according to tagger confidence**
Optimal confidence threshold must be determined. Sentence tagging accuracy of 98% possible when single tagging half of the tokens.

What are the supertags?

In CCG, supertags are the lexical entries.

But

$$\begin{array}{l}
 n_pp_c_obl_le \ \& \\
 \left[\begin{array}{l}
 \text{STEM} \quad \langle \text{"price"} \rangle, \\
 \text{SYNSEM} \quad \left[\begin{array}{l}
 \text{LKEYS} \quad \left[\begin{array}{l}
 \text{---COMPKEY} \quad _of_p_sel_rel \\
 \text{KEYREL.PRED} \quad _price_n_of_rel
 \end{array} \right] \\
 \text{PHON} \quad \left[\text{ONSET} \quad \text{con} \right]
 \end{array} \right]
 \end{array} \right]
 \end{array}$$

is not a good supertag.

We can use **letype**: `n_pp_c_obl_le`

But also:

subcat: `n_pp`

pos: `n`

What are the supertags?

	Coverage	Precision	F-score	Lexical Items	Lexical Ambiguity	Time (sec.)
Baseline	1.00	0.778	0.777	106.04	7.57	0.77

Table: Parser performance over *jhpstg-test*

What are the supertags?

	Coverage	Precision	F-score	Lexical Items	Lexical Ambiguity	Time (sec.)
Baseline	1.00	0.778	0.777	106.04	7.57	0.77
Gold letype	1.00	0.808	0.808	20.05	1.43	0.11
Gold subcat	1.00	0.805	0.804	28.46	2.03	0.13
Gold pos	1.00	0.800	0.795	55.56	3.97	0.26

Table: Parser performance over *jhpstg-test*

- Upper bound speed increase seven-fold
- Gold **subcat** is almost as fast as gold **letype**

What are the supertags?

	Coverage	Precision	F-score	Lexical Items	Lexical Ambiguity	Time (sec.)
Baseline	1.00	0.778	0.777	106.04	7.57	0.77
Gold letype	1.00	0.808	0.808	20.05	1.43	0.11
Single	0.67	0.765	0.550	19.36	1.38	0.09
Beta 0.001	0.71	0.767	0.589	33.16	2.37	0.11
Thresh 0.99	1.00	0.777	0.778	83.07	5.93	0.53

Table: Parser performance over *jhpstg-test*

What are the supertags?

	Coverage	Precision	F-score	Lexical Items	Lexical Ambiguity	Time (sec.)
Baseline	1.00	0.778	0.777	106.04	7.57	0.77
Gold subcat	1.00	0.805	0.804	28.46	2.03	0.13
Single	0.78	0.770	0.639	27.78	1.98	0.12
Beta 0.001	0.82	0.762	0.656	46.75	3.34	0.11
Thresh 0.99	1.00	0.777	0.777	84.73	6.05	0.56

Table: Parser performance over *jhpstg-test*

Thank You!