



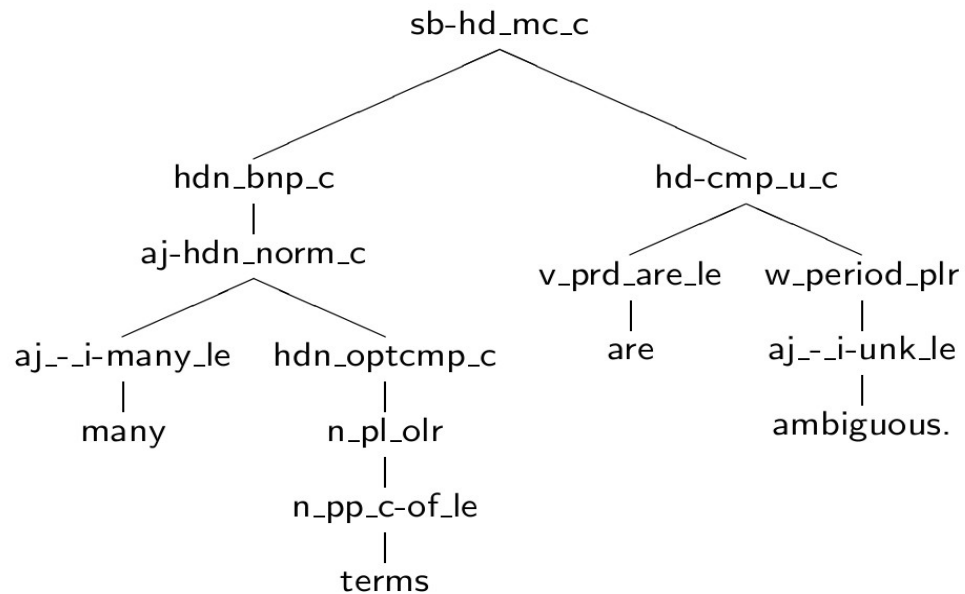
# Introducing CuteForce – Efficient HPSG Parsing

Gisle Ytrestøl – University of Oslo  
[gisley@ifi.uio.no](mailto:gisley@ifi.uio.no)



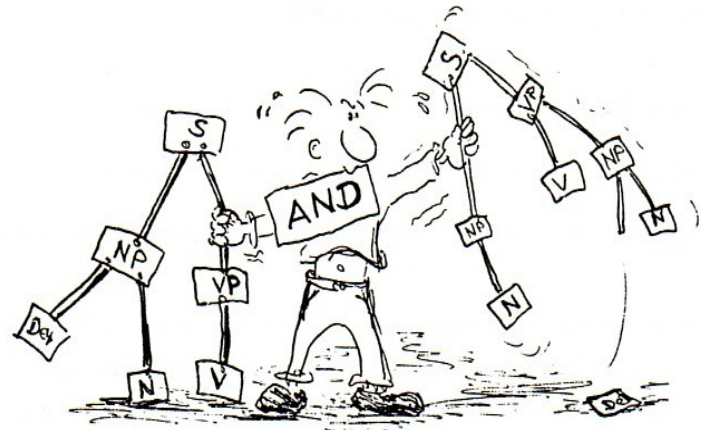
# Goal

- Produce derivations like this:



# Why? We already have PET/ERG!

- PET/ERG offers:
  - linguistically well-founded approach
  - coverage ~80 – 90% across domains: Wikipedia, GENIA, WSJ, et al
  - 500+ number of parses for long/medium length sentences
  - fully correct tree ranked best in ~55% of cases
  - an average ~5 sec wait for each sentence

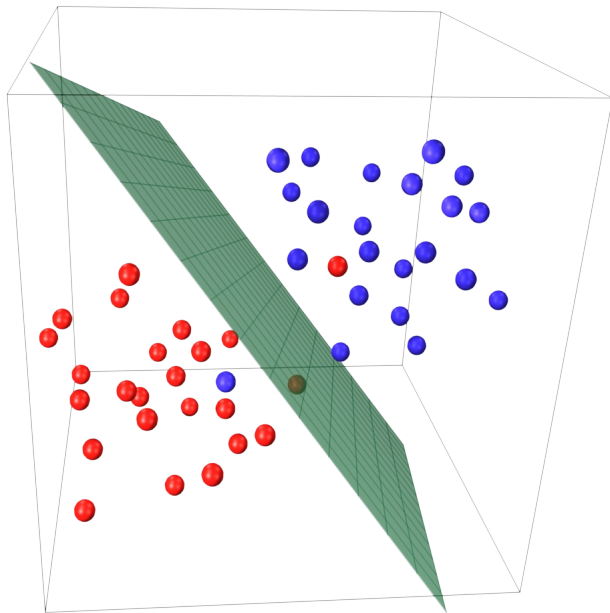


# CuteForce - A Cuter Approach

- Derives only one parse per sentence
- Knows very little about linguistics
- Fast (<5 ms per sentence – would theoretically parse Wikipedia in 120 cpu hours, not 14 years)
- Written in Java
- Light-weight



# CuteForce Parsing



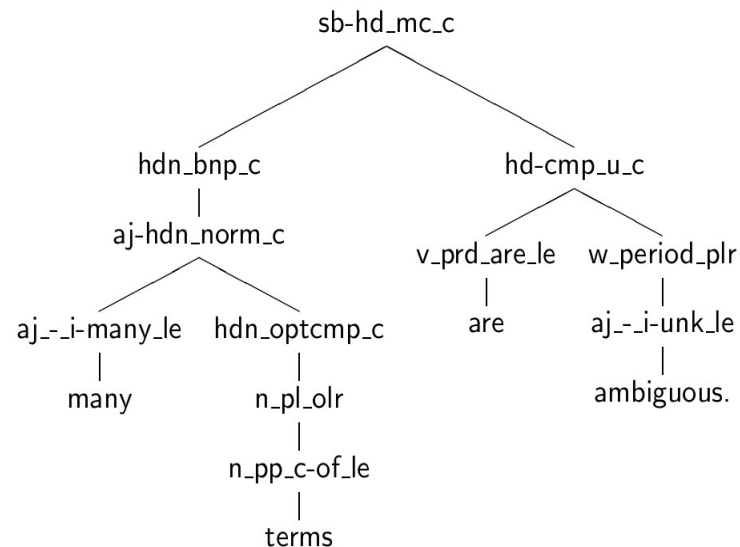
- Builds the HPSG representation incrementally in a shift-reduce fashion
- Inspired by the MaltParser (Nivre, 2007)
- Any HPSG structure can be built through CuteForce transitions
- Transitions are selected through SVM models

# CuteForce Parsing

- Learns this:

```
ARG #1=\\\"defined\\\" +PRED predsor +CLASS alpnae
#1 ]\"))))) (175324 hdn-aj_redrel_c 8.49123 15 20
nstruction_n1 -0.230751 15 16 (\"instructions\" 274
#1=\\\"instructions\\\" +PRED predsor +CLASS alpha
RM #1 ]\"))))) (175323 hd-cmp_u_c 6.12964 16 20 (666
\\\"91\\\" +ID *diff-list* +CARG #1=\\\"for\\\" +PRED
TRAIT native_trait +FORM #1 ]\")) (175322 hdn_bnp-v
.39161 17 20 (175316 v_prp_olr 0.299762 17 18 (675
\\\" +FROM \\\"95\\\" +ID *diff-list* +CARG #1=\\\"com
alized+lower ] +TRAIT native_trait +FORM #1 ]\"))
[ +TNT null_tnt +TO \\\"107\\\" +FROM \\\"106\\\" +ID
L - +CASE non_capitalized+lower ] +TRAIT native_tr
79986 19 20 (698 task_n1 0 19 20 (\"task\" 282 \"toke
=\\\"task\\\" +PRED predsor +CLASS alphabetic [ +IN
))))) (175930 hd-aj_scp-pr_c 39.3194 20 41
```

- Produces this:



# Data Structure

- Input buffer  $\beta$ :
  - Tuples of word forms/PTB tags/lexical types:
  - (many|JJ|aj\_-\_i-many\_le),  
(terms|NNS|n\_pp\_c-of\_le),  
(are|VBP|v\_prd\_are\_le),  
(ambiguous.|JJ|aj\_-\_i-nk\_le)
- Position index  $i$ :
  - The position index for  $\beta$ , starting at 1 for the first input triple.
- Active edges  $\alpha$ :
  - Stack of active edges
  - Binary Constituents whose right-branched daughter/mother is not yet assigned
  - Empty upon termination
- Passive edges  $\pi$ :
  - Stack of passive edges
  - Represents the partial HPSG structure during parsing, and the full HPSG representation upon termination



# Transition system

- ACTIVE (adds active edge to  $\alpha$ , adds lexical type  $\beta(i+1)$  to  $\pi$ , and increments  $i$ )
- PASSIVE (pops  $\alpha$  and adds binary passive edge to  $\pi$ )
- UNIT (adds unary passive edge to  $\pi$ )
- ACCEPT (Terminates the parse of the sentence.  $\pi$  represents the syntactic derivation of the sentence)



# Initial Parse State

•  $\pi$ :

aj\_-\_i-many\_le  
|  
many

- Input buffer  $\beta$ :
  - (many|JJ|aj\_-\_i-many\_le),  
(terms|NNS|n\_pp\_c-of\_le),  
(are|VBP|v\_prd\_are\_le),  
(ambiguous.|JJ|aj\_-\_i-nk\_le)
  - $i = 1$
- $\alpha = e$

# Active

• $\pi$ :

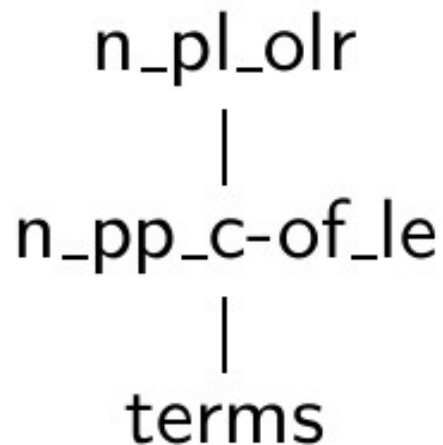
n\_pp\_c-of\_le  
|  
terms

- Input buffer  $\beta$ :
  - (many|JJ|aj\_-\_i-many\_le),  
(terms|NNS|n\_pp\_c-of\_le),  
(are|VBP|v\_prd\_are\_le),  
(ambiguous.|JJ|aj\_-\_i-nk\_le)
  - $i = 2$
- $\alpha$ :

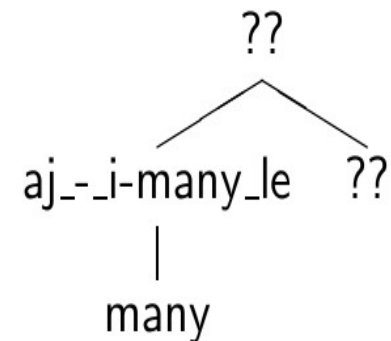
??  
/ \  
aj\_-\_i-many\_le ??  
|  
many

# Unit - (n\_pl\_olr)

• $\pi$ :



- Input buffer  $\beta$ :
  - (many|JJ|aj\_-\_i-many\_le),  
(terms|NNS|n\_pp\_c-of\_le),  
(are|VBP|v\_prd\_are\_le),  
(ambiguous.|JJ|aj\_-\_i-nk\_le)
  - $i = 2$
- $\alpha$ :



# Unit - (hdn\_optcmp\_c)

•  $\pi$ :

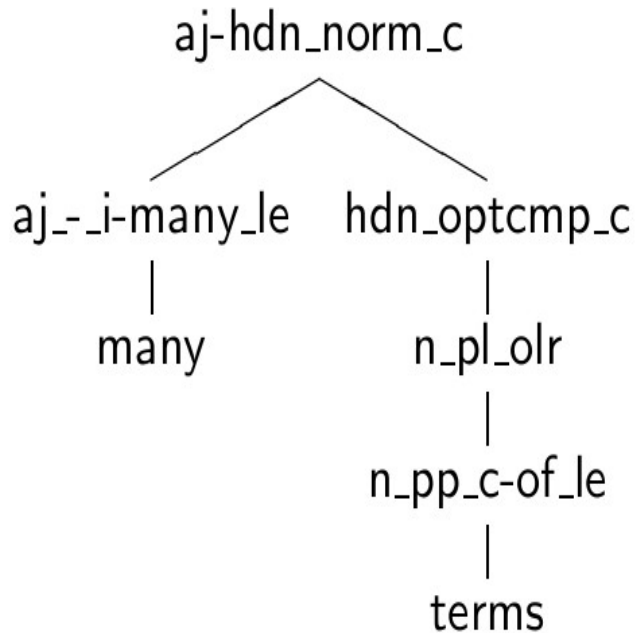
hdn\_optcmp\_c  
|  
n\_pl\_olr  
|  
n\_pp\_c-of\_le  
|  
terms

- Input buffer  $\beta$ :
  - (many|JJ|aj\_-\_i-many\_le),  
(terms|NNS|n\_pp\_c-of\_le),  
(are|VBP|v\_prd\_are\_le),  
(ambiguous.|JJ|aj\_-\_i-nk\_le)
  - $i = 2$
- $\alpha$ :

??  
/ \  
aj\_-\_i-many\_le ??  
|  
many

# Passive - (aj-hdn\_norm\_c)

• $\pi$ :



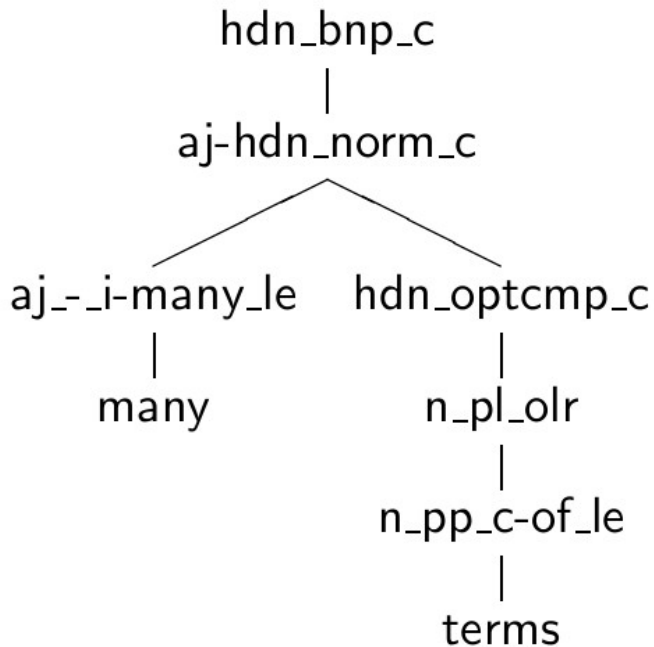
• Input buffer  $\beta$ :

- (many|JJ|aj\_-\_i-many\_le),  
(terms|NNS|n\_pp\_c-of\_le),  
(are|VBP|v\_prd\_are\_le),  
(ambiguous.|JJ|aj\_-\_i-nk\_le)
- $i = 2$

•  $\alpha$ :

# Unit - (hdn\_bnp\_c)

•  $\pi$ :



• Input buffer  $\beta$ :

- (many|JJ|aj\_-\_i-many\_le),  
(terms|NNS|n\_pp\_c-of\_le),  
(are|VBP|v\_prd\_are\_le),  
(ambiguous.|JJ|aj\_-\_i-nk\_le)
- $i = 2$

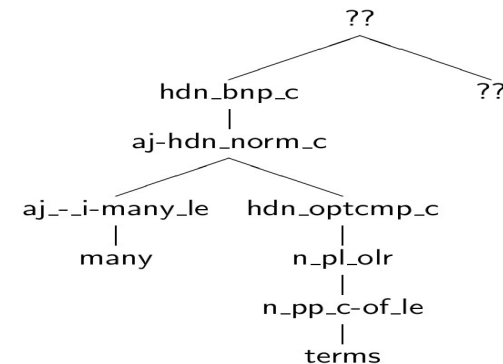
•  $\alpha$ :

# Active

## • $\pi$ :

v\_prd\_are\_le  
|  
are

- Input buffer  $\beta$ :
  - (many|JJ|aj\_-\_i-many\_le),  
(terms|NNS|n\_pp\_c-of\_le),  
(are|VBP|v\_prd\_are\_le),  
(ambiguous.|JJ|aj\_-\_i-nk\_le)
  - $i = 3$
- $\alpha$ :

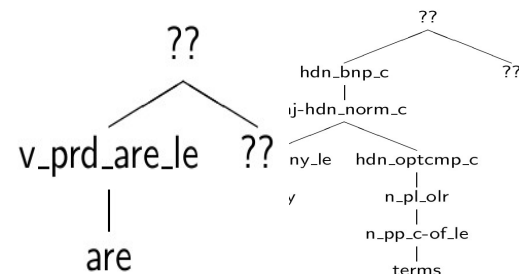


# Active

## • $\pi$ :

aj\_-\_i-unk\_le  
|  
ambiguous.

- Input buffer  $\beta$ :
  - (many|JJ|aj\_-\_i-many\_le),  
(terms|NNS|n\_pp\_c-of\_le),  
(are|VBP|v\_prd\_are\_le),  
(ambiguous.|JJ|aj\_-\_i-nk\_le)
  - $i = 4$
- $\alpha$ :



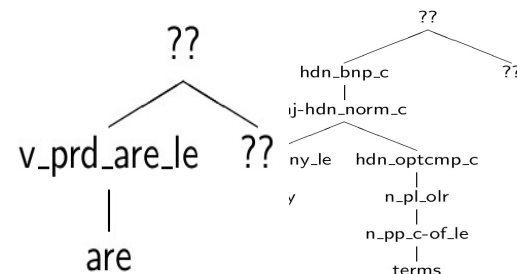


# Unit - (w\_period\_plr)

• $\pi$ :

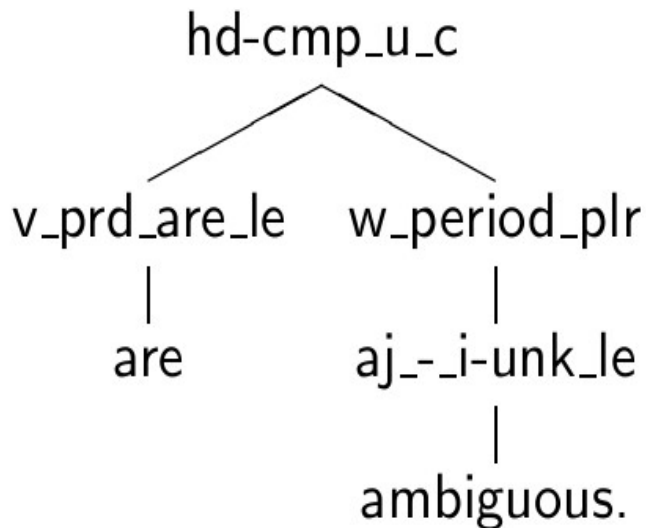
w\_period\_plr  
|  
aj\_-\_i-unk\_le  
|  
ambiguous.

- Input buffer  $\beta$ :
  - (many|JJ|aj\_-\_i-many\_le),  
(terms|NNS|n\_pp\_c-of\_le),  
(are|VBP|v\_prd\_are\_le),  
(ambiguous.|JJ|aj\_-\_i-nk\_le)
  - $i = 4$
- $\alpha$ :



# Passive - (hd-cmp\_u\_c)

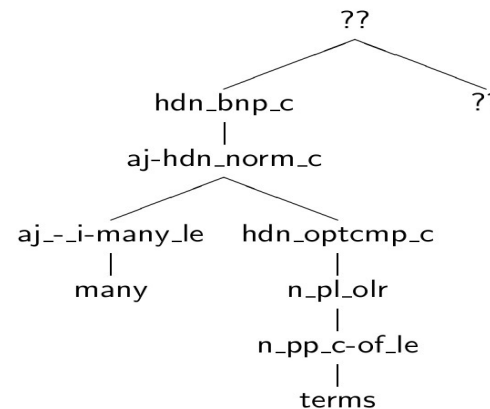
•  $\pi$ :



• Input buffer  $\beta$ :

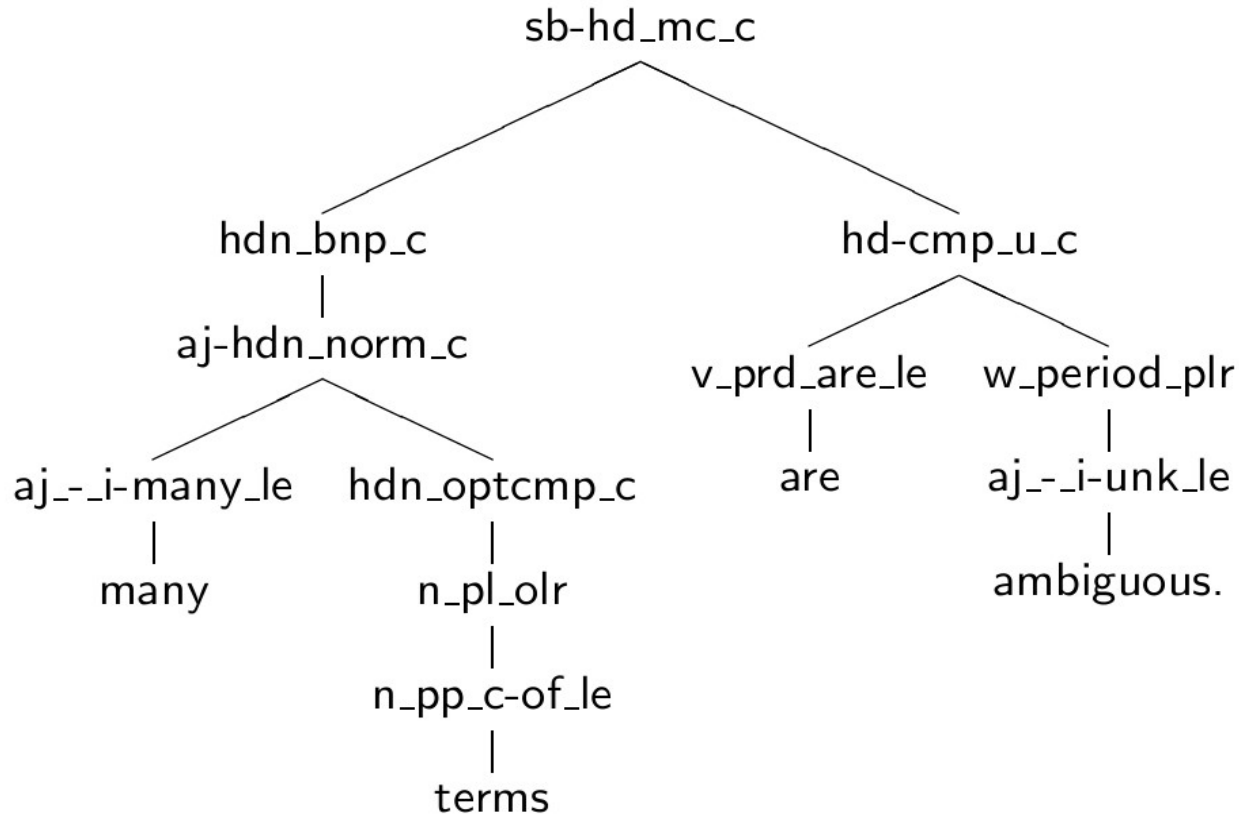
- (many|JJ|aj\_-\_i-many\_le),  
(terms|NNS|n\_pp\_c-of\_le),  
(are|VBP|v\_prd\_are\_le),  
(ambiguous.|JJ|aj\_-\_i-nk\_le)
- $i = 4$

•  $\alpha$ :

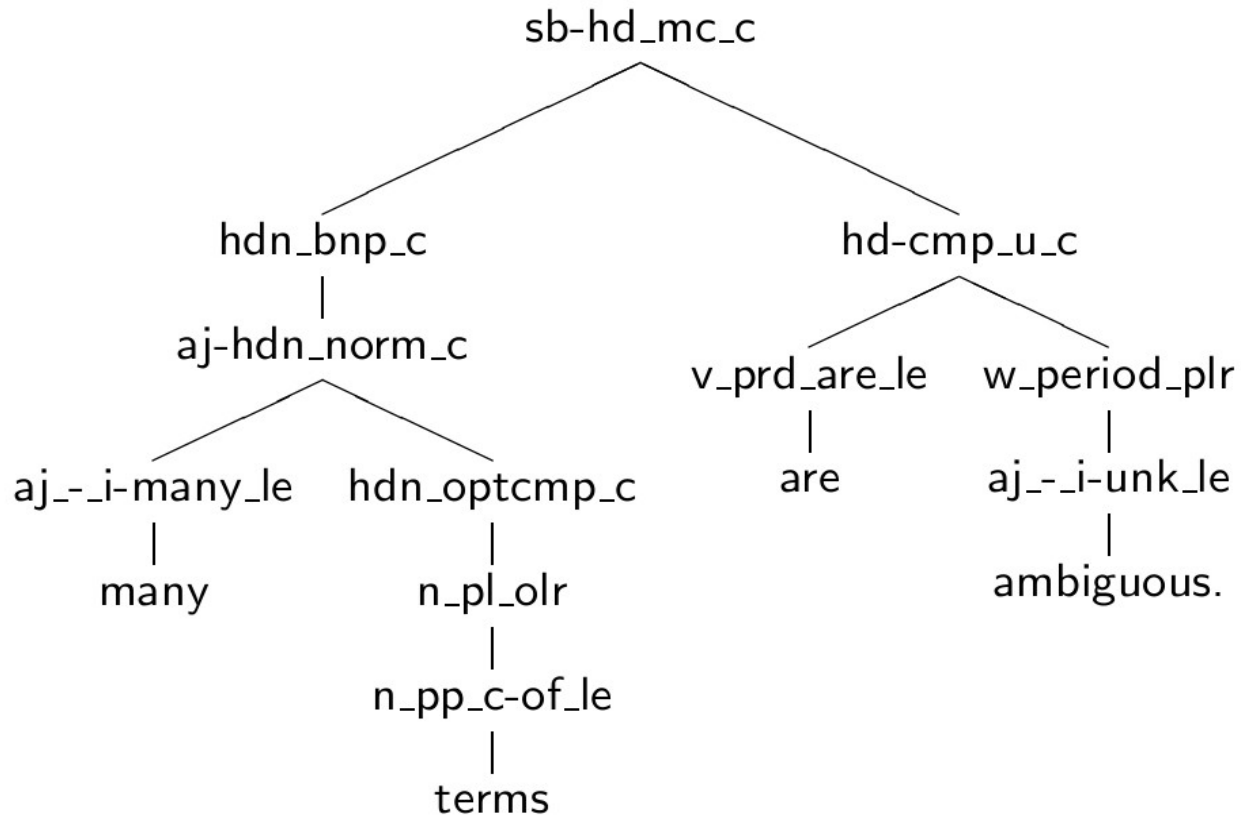


# Passive - (sb-hd\_mc\_c)

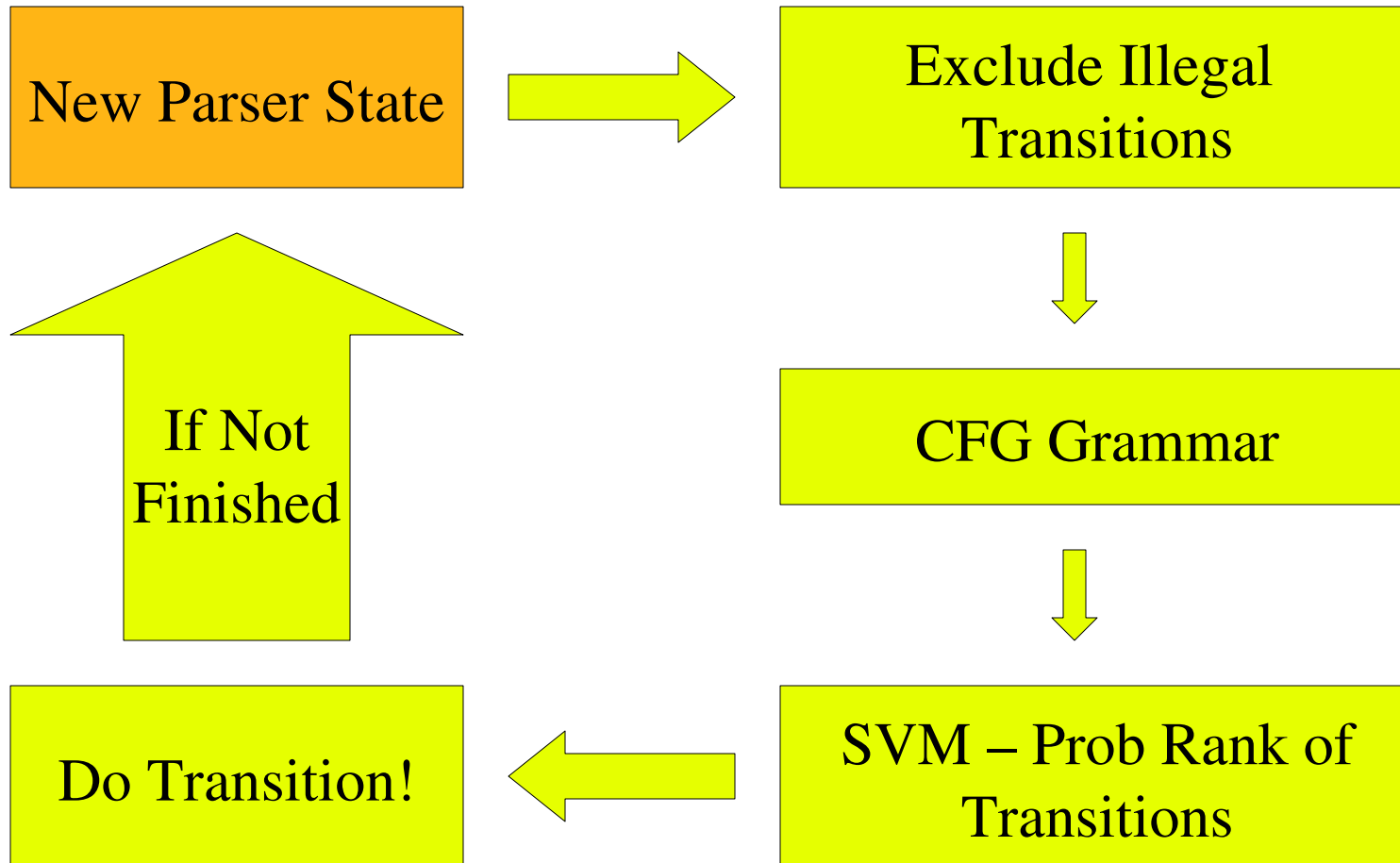
•π:



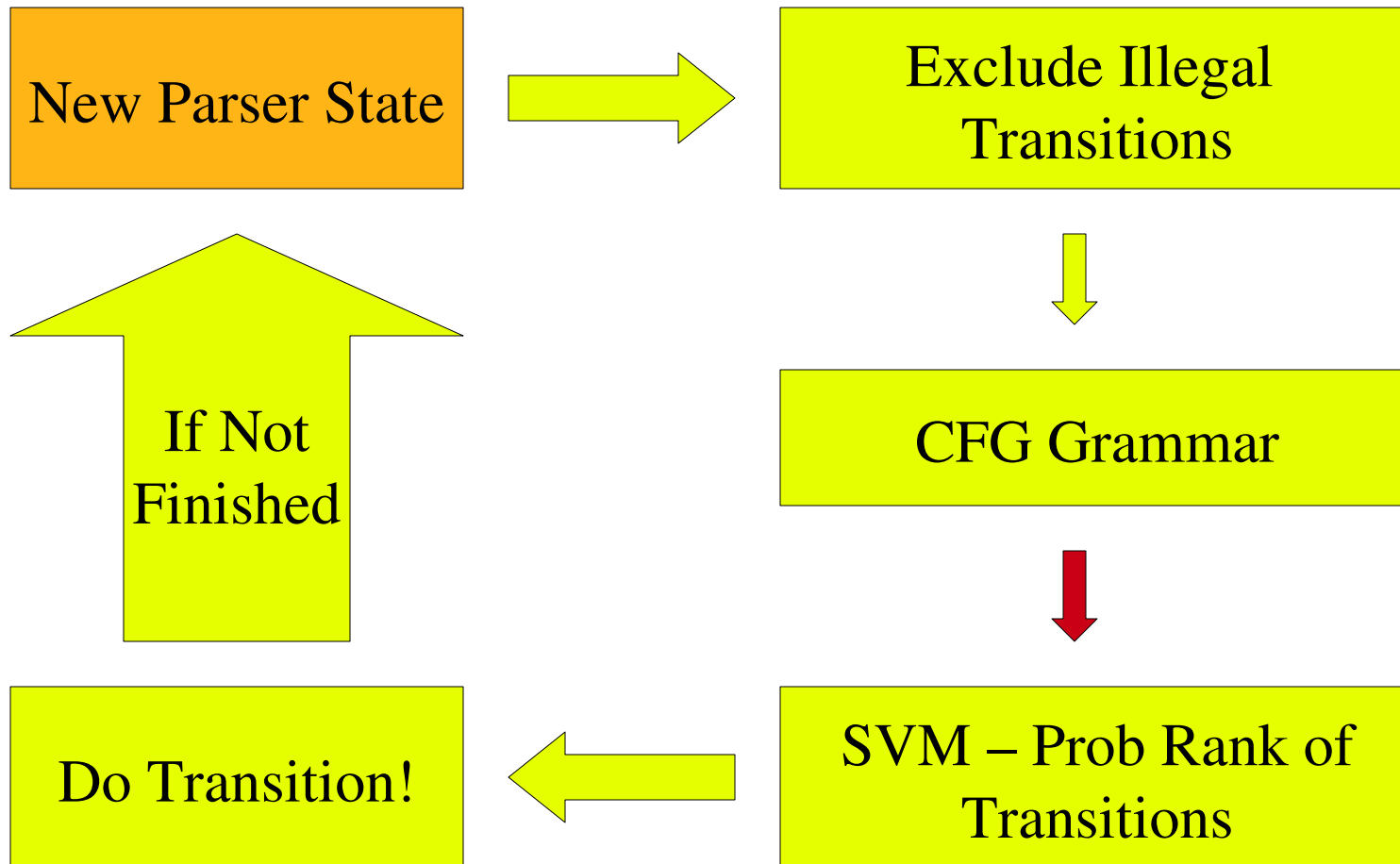
# Finished



# Oracle



# Common Point of Failure





# SuperTagging

- *Lexical Type: A feature structure containing category information for the lexical entry. Describes at least head type and valence information. (Dridan, 2009)*
- Input buffer:
  - Word form (ERG Gold Standard Tokenization)
  - PTB Tags (Training/Test data tagged by TnT)
  - Lexical Types (Gold Standard for Training, Test data tagged by C&C SuperTagger)

# SuperTagging

- *<part-of-speech>\_<subcategorisation>\_<description>\_le*
- Current C&C Accuracy:
  - LT Accuracy: 0.86
  - POS & SUB Accuracy: 0.93
  - POS Accuracy: 0.97







# Preliminary Results\* / Supertagged Input

- |                           |                              |
|---------------------------|------------------------------|
| • <u>CFG Constraints:</u> | • <u>No CFG Constraints:</u> |
| • Coverage: ~0.71         | • Coverage: ~0.99            |
| • Exact Match: ~0.39      | • Exact Match: ~0.26         |
| • Parseval f-score: ~0.82 | • Parseval f-score: ~0.76    |
| • Sentence Length:        | • Sentence Length:           |
| – Avg # of tokens: ~ 14.0 | – Avg # of tokens: ~ 14.0    |
| – Parsed: ~ 11.2          | – Parsed: ~ 14.0             |
| – Failed: ~ 20.9          | – Failed: ~ 18.0             |

\*Lexical types not evaluated

# Preliminary Results\* / Gold Standard Buffer Input

- CFG Constraints:
  - Coverage: ~0.72
  - Exact Match: ~0.50
  - Parseval f-score: ~0.89
  - Sentence Length:
    - Avg # of tokens: ~ 14.0
    - Parsed: ~ 11.2
    - Failed: ~ 21.3
- No CFG Constraints:
  - Coverage: ~0.99
  - Exact Match: ~0.35
  - Parseval f-score: ~0.84
  - Sentence Length:
    - Avg # of tokens: ~ 14.0
    - Parsed: ~ 14.0
    - Failed: ~ 17.8

**\*Lexical types not evaluated**

# Further Ways of Improving CF

- Sanity testing/bug fixing
- Deterministic approaches (cheap):
  - Improve feature model
  - Further experiments on partitioning of feature model
- Non-deterministic approach (expensive):
  - Back-tracking
- Unification validation using Pet (expensive)





# References

- Rebecca Dridan. Using lexical statistics to improve HPSG parsing. PhD thesis, Saarland University, 2009
- Joakim Nivre, Jens Nilsson, Johan Hall, Atanas Chanev, Gulsen Eryigit, Sandra Kubler, Svetoslav Marinov, and Erwin Marsi. Maltparser: A language-independent system for data-driven dependency parsing. *Natural Language Engineering*, 13(1):1-41, 2007.