

## Training Data Selection

- **Basic Filtering:** Keep balanced lines of  $\leq 95$  words with all words  $\leq 25$  characters
- **Giga-FrEn Filtering:** Keep balanced lines of  $\leq 50$  words containing a French word seen  $< 20$  times
- Number of parallel sentences:

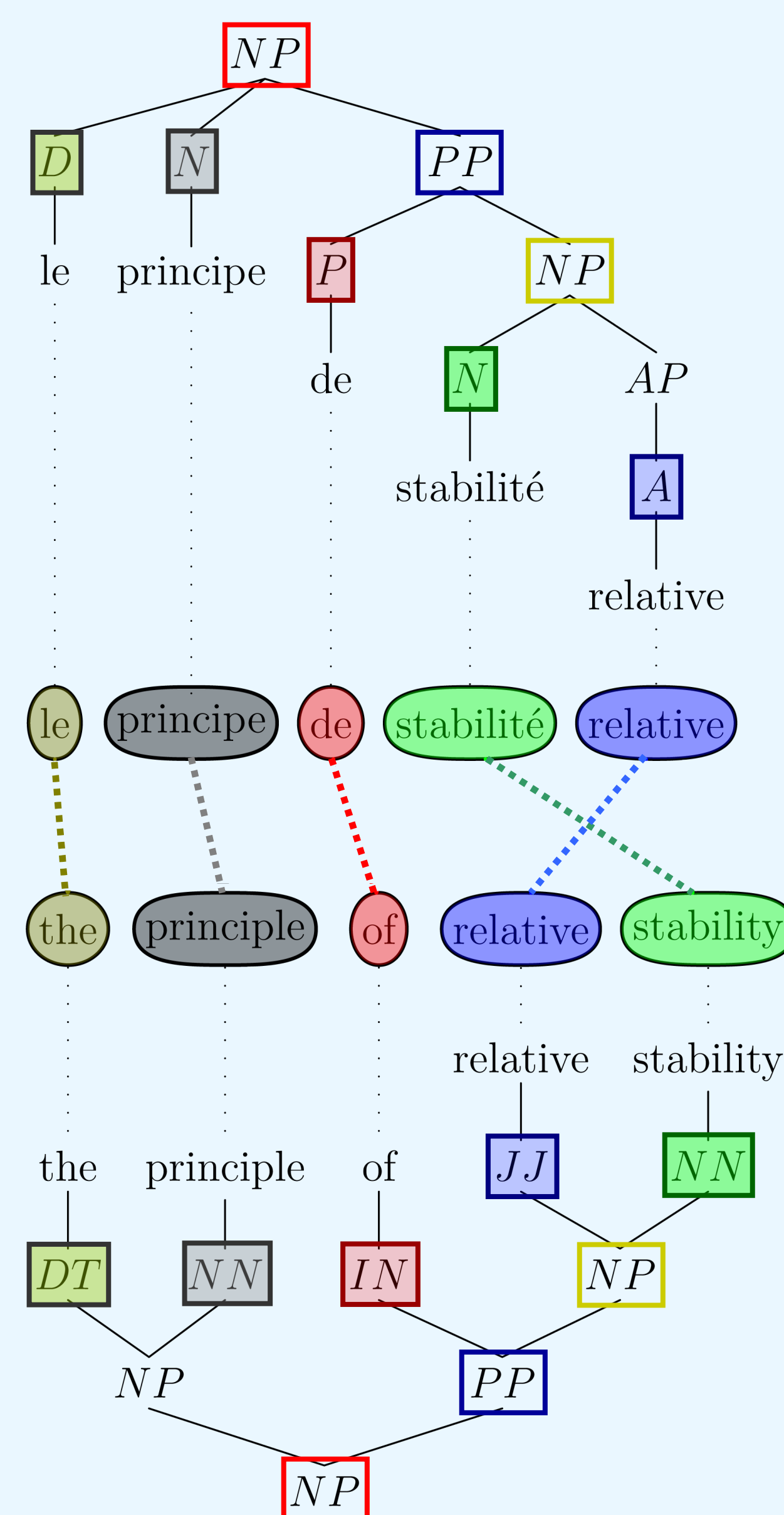
Corpus	Released	Used
Europarl	1,825,077	1,614,111
News Comm	115,562	95,138
UN Docs	12,317,600	9,352,232
Giga-FrEn	22,520,400	2,839,466
<b>Total</b>	<b>38,778,639</b>	<b>13,900,947</b>

- **+1.2 to +1.9 BLEU** and **-45% OOV rate** over using WMT 2010 data only

## External Tools

- MGIZA++
- Stanford tokenizer
- Berkeley parser
- Joshua decoder
- SRI LM toolkit
- ZMERT

## Grammar Extraction



- Extract SCFG rules as well as standard SMT phrase pairs
- When duplicates occur, keep syntactic version but share counts
- Give SMT phrase pairs dummy syntactic label

	le	principe	de	stabilité	relative
the	■				
principe		■			
of			■		
relative				■	
stability					■

$N::NN \rightarrow [\text{stabilité}]::[\text{stability}]$   
 $NP::NP \rightarrow [\text{stabilité relative}]::[\text{relative stability}]$   
 $NP::NP \rightarrow [N^1 A^2]::[JJ^2 NN^1]$   
 $PHR::PHR \rightarrow [\text{principe de}]::[\text{principe of}]$

## Grammar Filtering

- Always keep all phrase pairs matching test set
- **Method 1:** Keep the 10,000 most frequently extracted hierarchical rules of all types
- **Method 2:** Keep only the 2000 most frequent abstract rules, but the 100,000 most frequent partially lexicalized rules

## Results

- Score differences between grammar filtering techniques inconclusive (**-0.4 to +0.8 BLEU**)
- Method 1 uses more rule instances; Method 2 uses more unique rules
- Automatic metric scores on dev test sets:

System	newstest-2009		newstest-2010	
	METR	BLEU	METR	BLEU
WMT '10 Method 1	54.94	24.77	56.66	25.78
WMT '10 Method 2	55.16	24.88	56.89	26.05
WMT '11 Method 1	55.82	26.02	58.13	27.71
WMT '11 Method 2	55.77	26.01	57.88	27.38