

# **A General-Purpose Rule Extractor for SCFG-Based Machine Translation**

**Greg Hanneman, Michelle Burroughs, and Alon Lavie**

Language Technologies Institute  
Carnegie Mellon University

Fifth Workshop on Syntax and Structure in Statistical Translation  
June 23, 2011



**Carnegie Mellon**

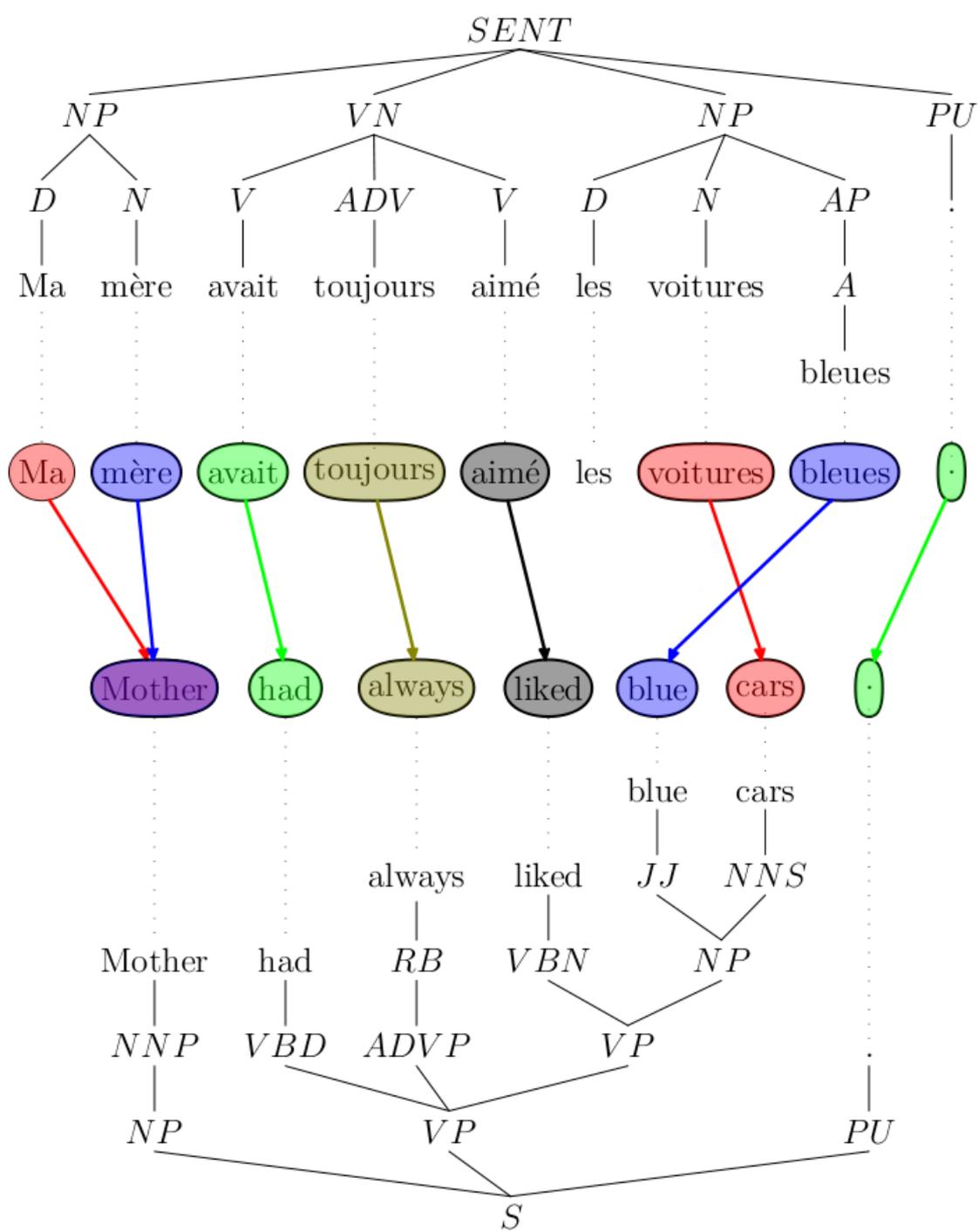
# SCFG Grammar Extraction

- Inputs:
  - Word-aligned sentence pair
  - Constituency parse trees on one or both sides
- Outputs:
  - Set of SCFG rules derivable from the inputs, possibly according to some constraints
- Implemented by:
  - Hiero [Chiang 2005]
  - Chiang [2010]
  - SAMT [Zollmann and Venugopal 2006]
  - GHKM [Galley et al. 2004]
  - Stat-XFER [Lavie et al. 2008]

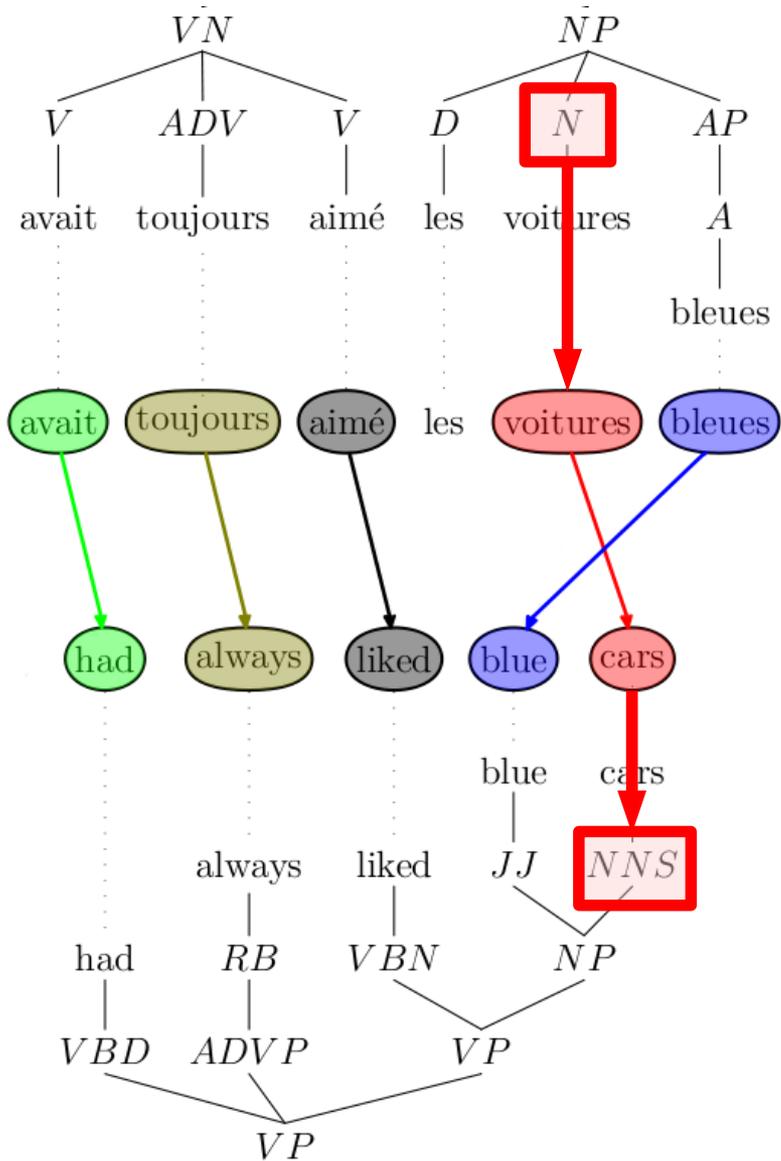
# SCFG Grammar Extraction

- Our goals:
  - Support for two parse trees by default
  - Extract greatest number of syntactic rules...
  - Without violating constituent boundaries
- Achieved with:
  - Multiple node alignments
  - Virtual nodes
  - Multiple right-hand-side decompositions

**First grammar extractor to do all three**

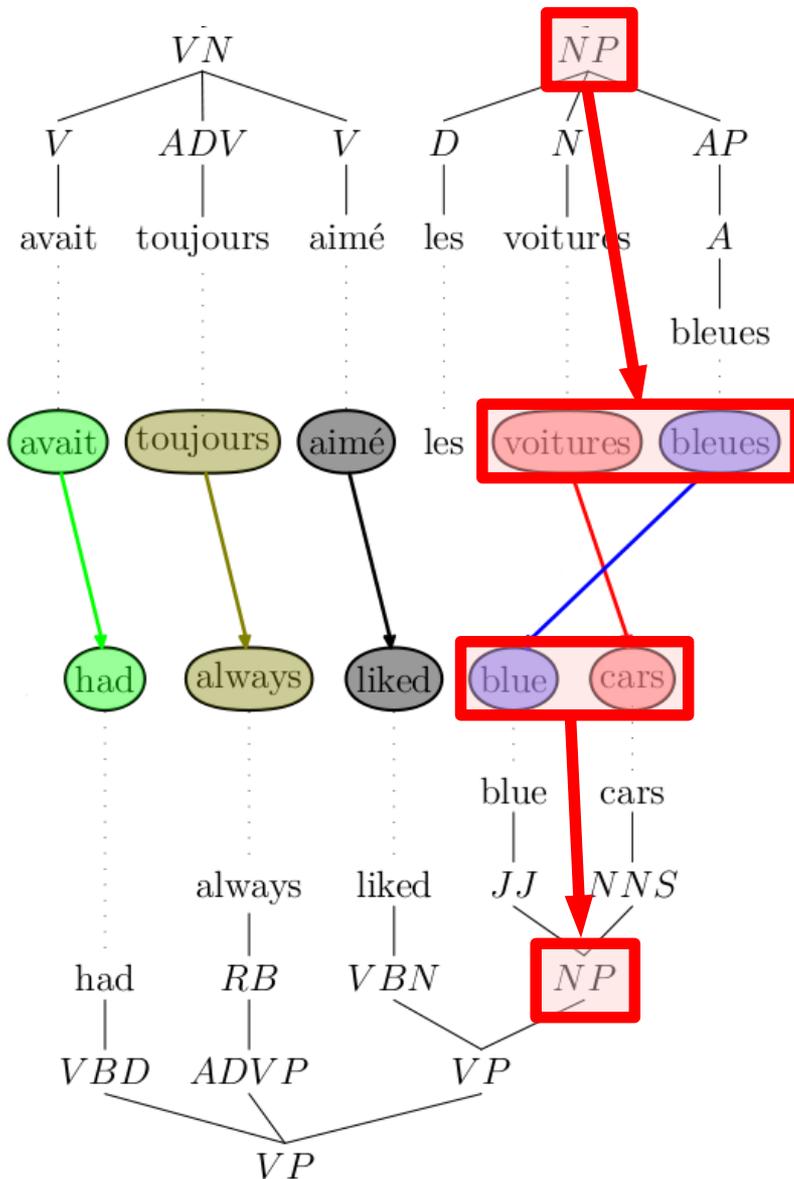


# Basic Node Alignment



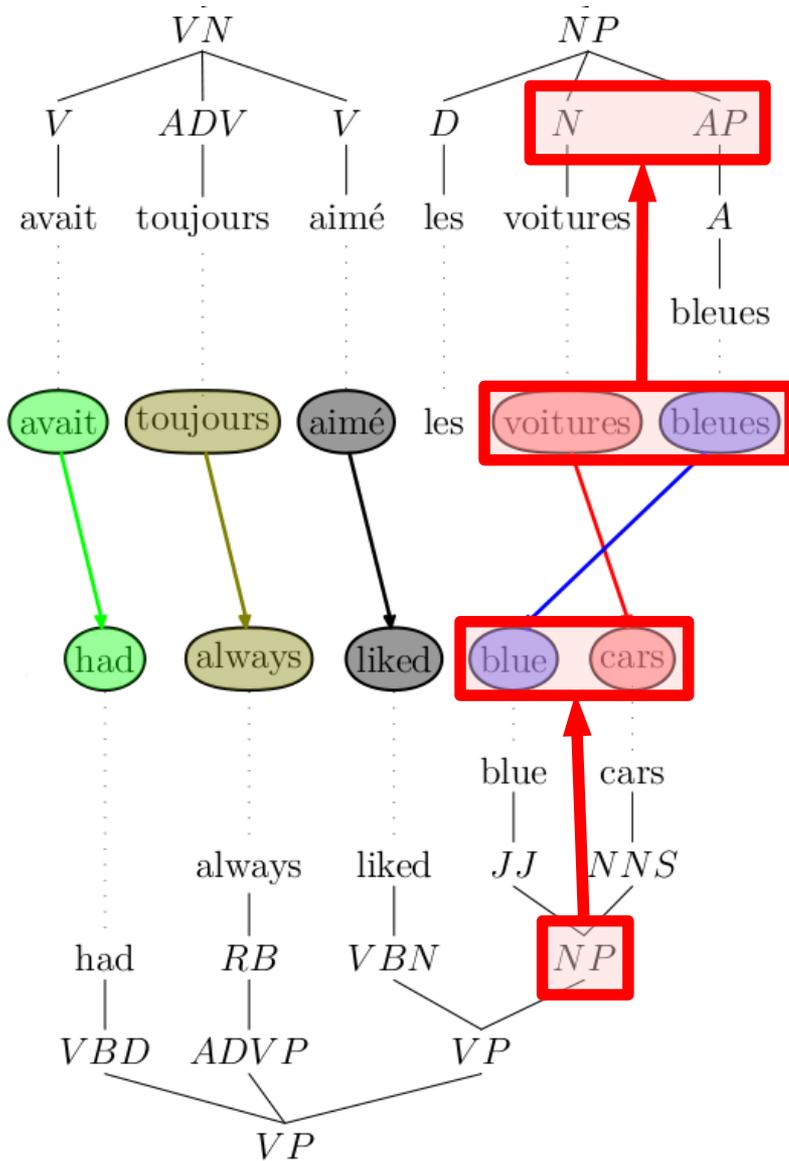
- Word alignment consistency constraint from phrase-based SMT

# Basic Node Alignment



- Word alignment consistency constraint from phrase-based SMT

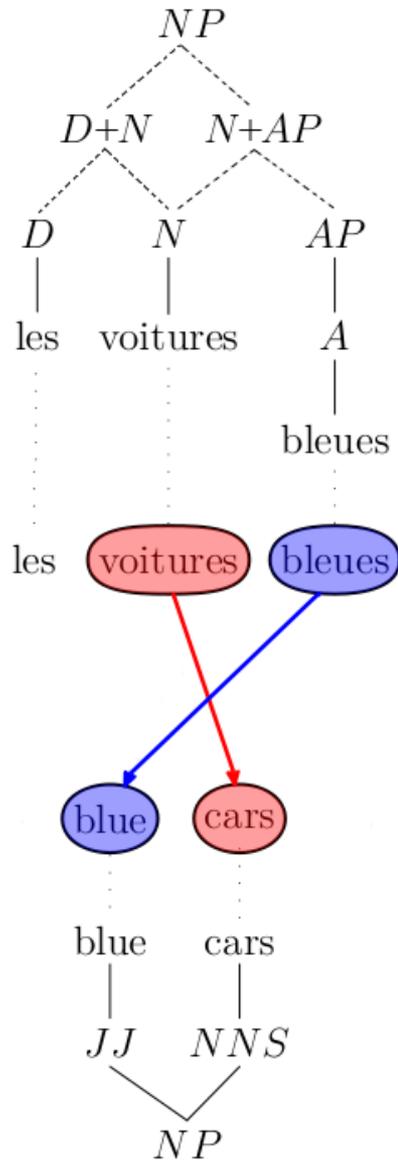
# Virtual Nodes



- Consistently aligned consecutive children of the same parent

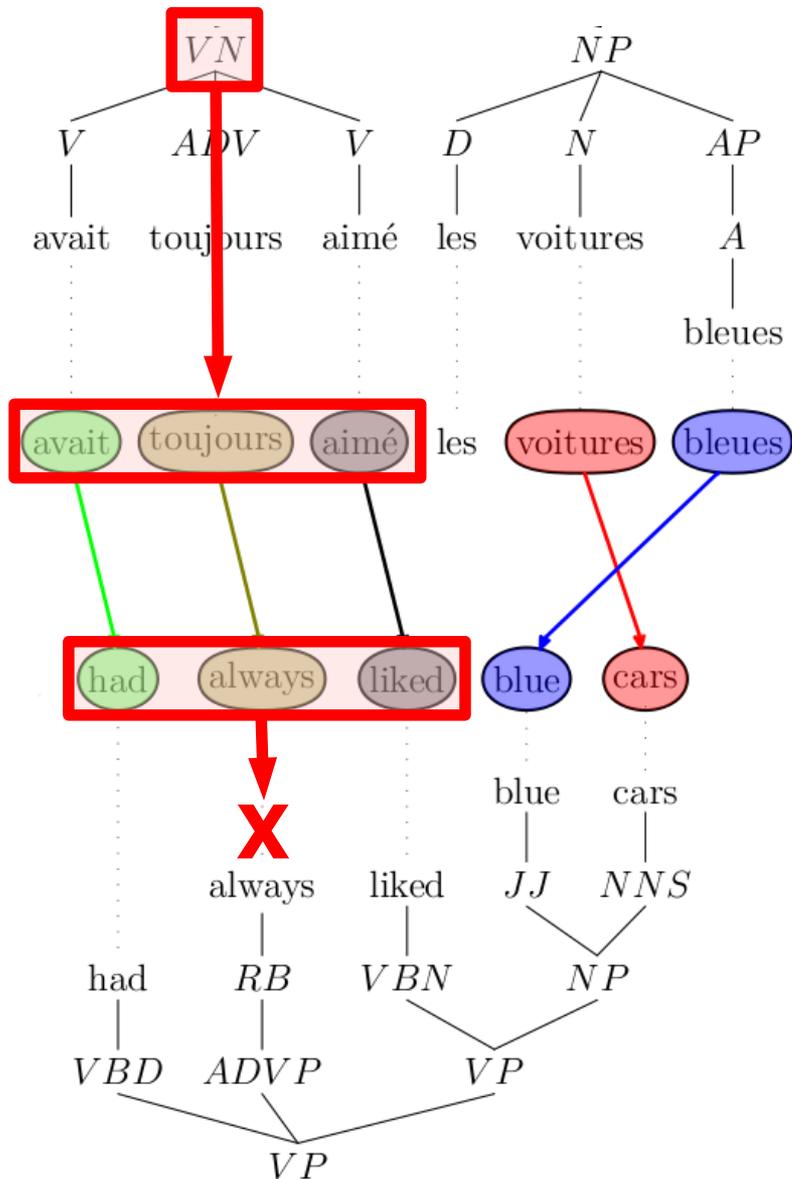


# Virtual Nodes



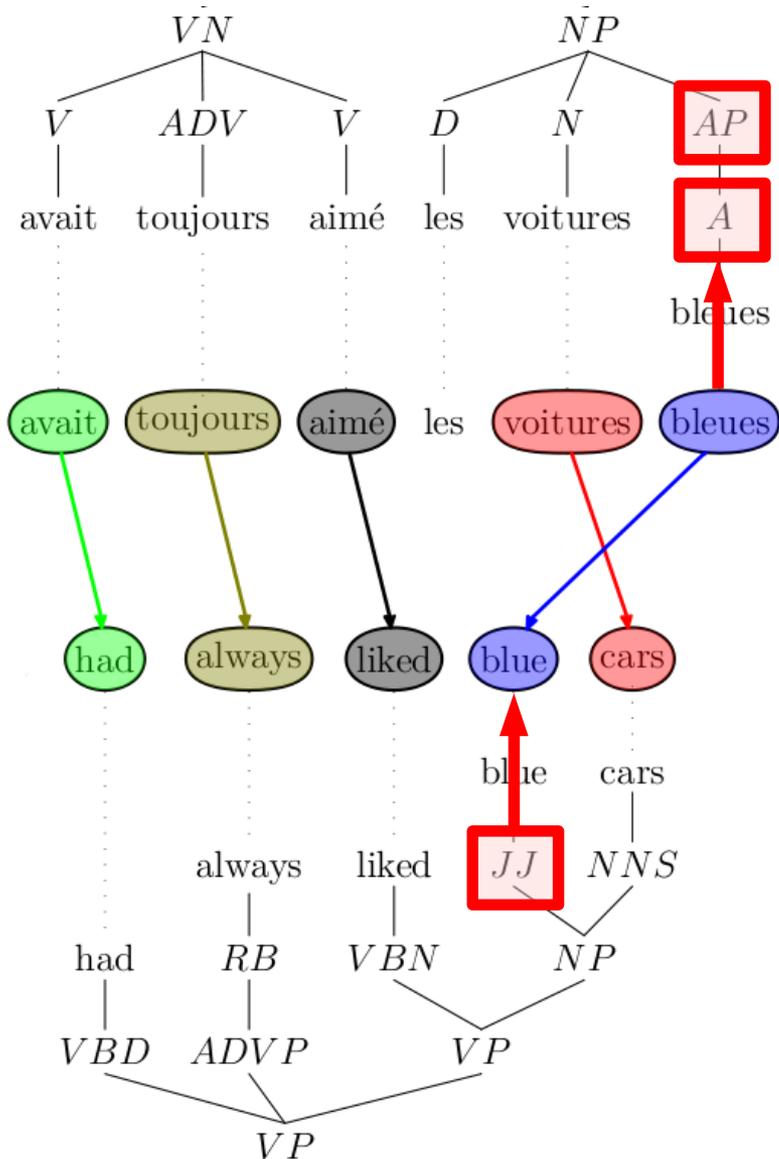
- Consistently aligned consecutive children of the same parent
- New intermediate node inserted in tree
- Virtual nodes may overlap
- Virtual nodes may align to any type of node

# Syntax Constraints



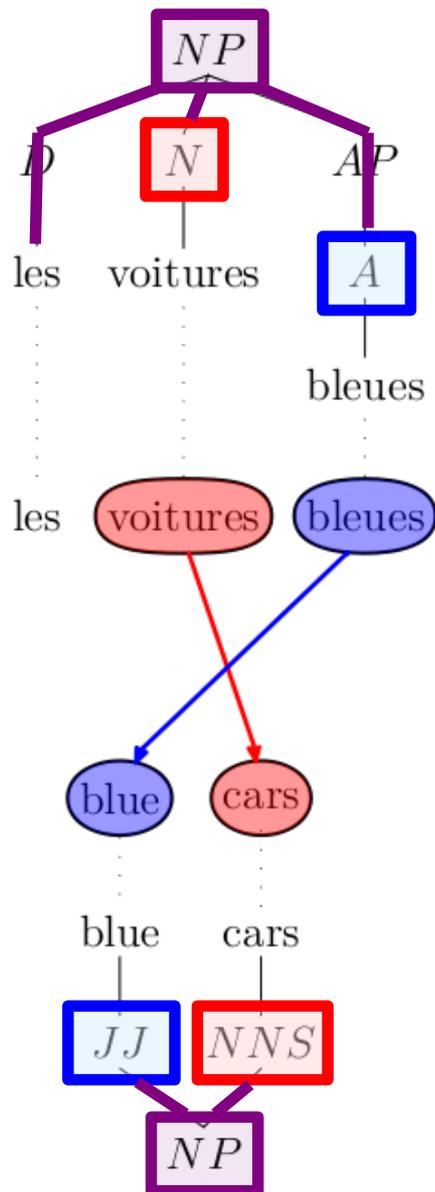
- Consistent word alignments  $\neq$  node alignment
- Virtual nodes may not cross constituent boundaries

# Multiple Alignment



- Nodes with multiple consistent alignments keep all of them

# Basic Grammar Extraction



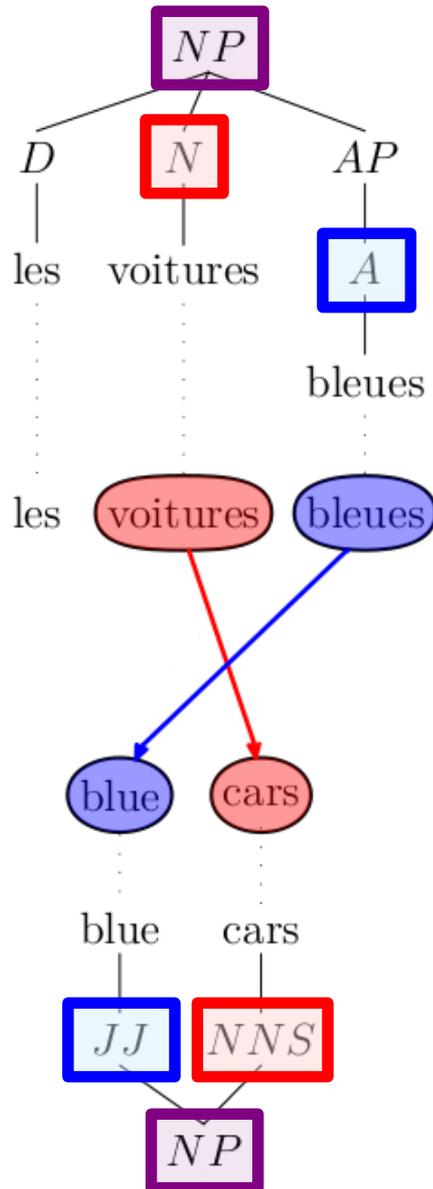
- Aligned node pair is LHS; aligned subnodes are RHS

$NP::NP \rightarrow [les\ N^1\ A^2]::[JJ^2\ NNS^1]$

$N::NNS \rightarrow [voitures]::[cars]$

$A::JJ \rightarrow [bleues]::[blue]$

# Multiple Decompositions



- All possible right-hand sides are extracted

$NP::NP \rightarrow [les N^1 A^2]::[JJ^2 NNS^1]$

$NP::NP \rightarrow [les N^1 bleues]::[blue NNS^1]$

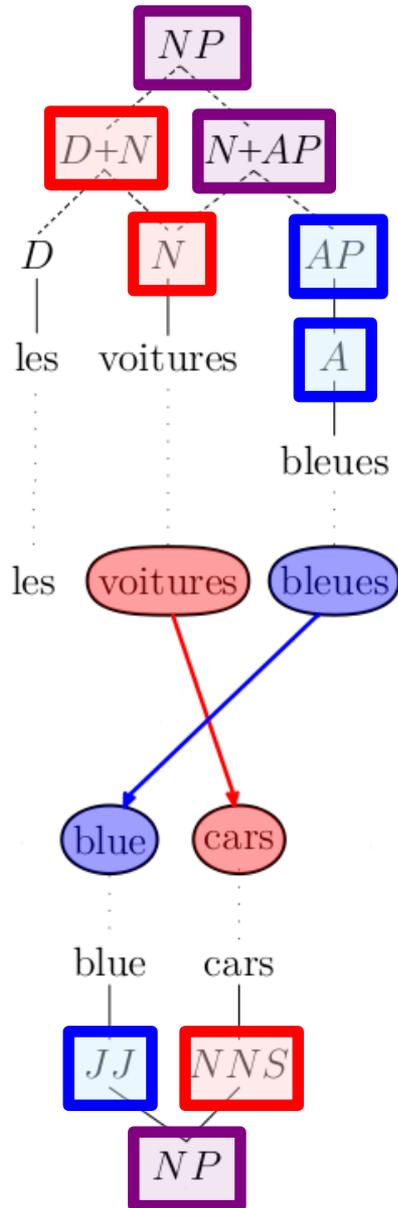
$NP::NP \rightarrow [les voitures A^2]::[JJ^2 cars]$

$NP::NP \rightarrow [les voitures bleues]::[blue cars]$

$N::NNS \rightarrow [voitures]::[cars]$

$A::JJ \rightarrow [bleues]::[blue]$

# Multiple Decompositions



- NP::NP  $\rightarrow$  [les N+AP<sup>1</sup>]::[NP<sup>1</sup>]
- NP::NP  $\rightarrow$  [D+N<sup>1</sup> AP<sup>2</sup>]::[JJ<sup>2</sup> NNS<sup>1</sup>]
- NP::NP  $\rightarrow$  [D+N<sup>1</sup> A<sup>2</sup>]::[JJ<sup>2</sup> NNS<sup>1</sup>]
- NP::NP  $\rightarrow$  [les N<sup>1</sup> AP<sup>2</sup>]::[JJ<sup>2</sup> NNS<sup>1</sup>]
- NP::NP  $\rightarrow$  [les N<sup>1</sup> A<sup>2</sup>]::[JJ<sup>2</sup> NNS<sup>1</sup>]
- NP::NP  $\rightarrow$  [D+N<sup>1</sup> bleues]::[blue NNS<sup>1</sup>]
- NP::NP  $\rightarrow$  [les N<sup>1</sup> bleues]::[blue NNS<sup>1</sup>]
- NP::NP  $\rightarrow$  [les voitures AP<sup>2</sup>]::[JJ<sup>2</sup> cars]
- NP::NP  $\rightarrow$  [les voitures A<sup>2</sup>]::[JJ<sup>2</sup> cars]
- NP::NP  $\rightarrow$  [les voitures bleues]::[blue cars]
- D+N::NNS  $\rightarrow$  [les N<sup>1</sup>]::[NNS<sup>1</sup>]
- D+N::NNS  $\rightarrow$  [les voitures]::[cars]
- N+AP::NP  $\rightarrow$  [N<sup>1</sup> AP<sup>2</sup>]::[JJ<sup>2</sup> NNS<sup>1</sup>]
- N+AP::NP  $\rightarrow$  [N<sup>1</sup> A<sup>2</sup>]::[JJ<sup>2</sup> NNS<sup>1</sup>]
- N+AP::NP  $\rightarrow$  [N<sup>1</sup> bleues]::[blue NNS<sup>1</sup>]
- N+AP::NP  $\rightarrow$  [voitures AP<sup>2</sup>]::[JJ<sup>2</sup> cars]
- N+AP::NP  $\rightarrow$  [voitures A<sup>2</sup>]::[JJ<sup>2</sup> cars]
- N+AP::NP  $\rightarrow$  [voitures bleues]::[blue cars]
- N::NNS  $\rightarrow$  [voitures]::[cars]
- AP::JJ  $\rightarrow$  [A<sup>1</sup>]::[JJ<sup>1</sup>]
- AP::JJ  $\rightarrow$  [bleues]::[blue]
- A::JJ  $\rightarrow$  [bleues]::[blue]

# Constraints

- Max rank of phrase pair rules
- Max rank of hierarchical rules
- Max number of siblings in a virtual node
- Whether to allow unary chain rules

$$\text{NP}::\text{NP} \rightarrow [\text{PRO}^1]::[\text{PRP}^1]$$

- Whether to allow “triangle” rules

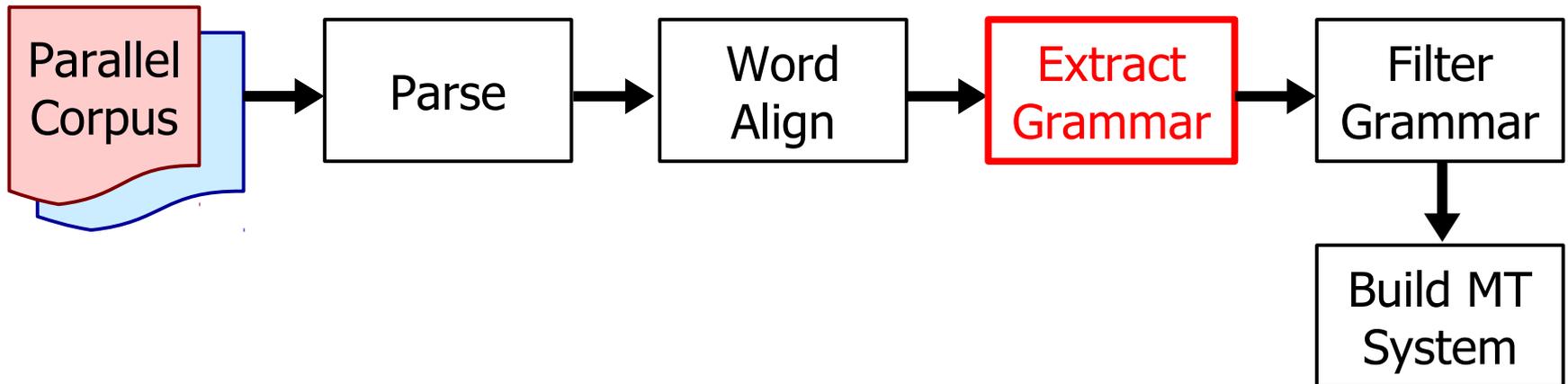
$$\text{AP}::\text{JJ} \rightarrow [\text{A}^1]::[\text{JJ}^1]$$

# Comparison to Related Work

	<b>Tree Constr.</b>	<b>Multiple Aligns</b>	<b>Virtual Nodes</b>	<b>Multiple Decomp.</b>
<b>Hiero</b>	No	—	—	Yes
<b>Stat-XFER</b>	Yes	No	Some	No
<b>GHKM</b>	Yes	No	No	Yes
<b>SAMT</b>	No	No	Yes	Yes
<b>Chiang [2010]</b>	No	No	Yes	Yes
<b>This work</b>	Yes	Yes	Yes	Yes

# Experimental Setup

- Train: FBIS Chinese–English corpus
- Tune: NIST MT 2006
- Test: NIST MT 2003



# Extraction Configurations

- **Baseline:**
  - Stat-XFER exact tree-to-tree extractor
  - Single decomposition with minimal rules
- **Multi:**
  - Add multiple alignments and decompositions
- **Virt short:**
  - Add virtual nodes; max rule length 5
- **Virt long:**
  - Max rule length 7

# Number of Rules Extracted

	Tokens		Types	
	Phrase	Hierarc.	Phrase	Hierarc.
<b>Baseline</b>	6,646,791	1,876,384	1,929,641	767,573
<b>Multi</b>	8,709,589	6,657,590	2,016,227	3,590,184
<b>Virt short</b>	10,190,487	14,190,066	2,877,650	8,313,690
<b>Virt long</b>	10,288,731	22,479,863	2,970,403	15,750,695

# Number of Rules Extracted

	Tokens		Types	
	Phrase	Hierarc.	Phrase	Hierarc.
<b>Baseline</b>	6,646,791	1,876,384	1,929,641	767,573
<b>Multi</b>	8,709,589	6,657,590	2,016,227	3,590,184
<b>Virt short</b>	10,190,487	14,190,066	2,877,650	8,313,690
<b>Virt long</b>	10,288,731	22,479,863	2,970,403	15,750,695

- Multiple alignments and decompositions:
  - Four times as many hierarchical rules
  - Small increase in number of phrase pairs

# Number of Rules Extracted

	Tokens		Types	
	Phrase	Hierarc.	Phrase	Hierarc.
<b>Baseline</b>	6,646,791	1,876,384	1,929,641	767,573
<b>Multi</b>	8,709,589	6,657,590	2,016,227	3,590,184
<b>Virt short</b>	10,190,487	14,190,066	2,877,650	8,313,690
<b>Virt long</b>	10,288,731	22,479,863	2,970,403	15,750,695

- Multiple decomp and virtual nodes:
  - 20 times as many hierarchical rules
  - Stronger effect on phrase pairs
  - 46% of rule types use virtual nodes

# Number of Rules Extracted

	Tokens		Types	
	Phrase	Hierarc.	Phrase	Hierarc.
<b>Baseline</b>	6,646,791	1,876,384	1,929,641	767,573
<b>Multi</b>	8,709,589	6,657,590	2,016,227	3,590,184
<b>Virt short</b>	10,190,487	14,190,066	2,877,650	8,313,690
<b>Virt long</b>	10,288,731	22,479,863	2,970,403	15,750,695

- Proportion of singletons mostly unchanged
- Average hierarchical rule count drops

# Rule Filtering for Decoding

- All phrase pair rules that match test set
- Most frequent hierarchical rules:
  - Top 10,000 of all types
  - Top 100,000 of all types
  - Top 5,000 fully abstract + top 100,000 partially lexicalized

$VP::ADJP \rightarrow [VV^1 VV^2]::[RB^1 VBN^2]$

$NP::NP \rightarrow [2000 \text{ 年 } NN^1]::[the 2000 NN^1]$

# Results: Metric Scores

- NIST MT 2003 test set

<b>System</b>	<b>Filter</b>	<b>BLEU</b>	<b>METR</b>	<b>TER</b>
<b>Baseline</b>	<b>10k</b>	24.39	54.35	68.01
<b>Multi</b>	<b>10k</b>	24.28	53.58	65.30
<b>Virt short</b>	<b>10k</b>	25.16	54.33	66.25
<b>Virt long</b>	<b>10k</b>	25.74	54.55	65.52

- Strict grammar filtering: extra phrase pairs help improve scores

# Results: Metric Scores

- NIST MT 2003 test set

<b>System</b>	<b>Filter</b>	<b>BLEU</b>	<b>METR</b>	<b>TER</b>
<b>Baseline</b>	<b>5k+100k</b>	25.95	54.77	66.27
<b>Virt short</b>	<b>5k+100k</b>	26.08	54.58	64.32
<b>Virt long</b>	<b>5k+100k</b>	25.83	54.35	64.55

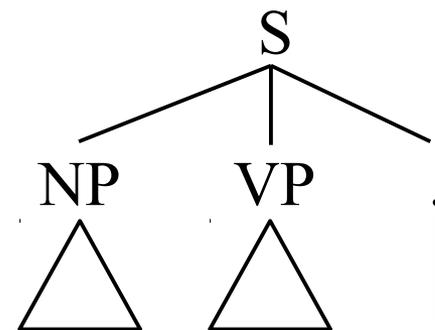
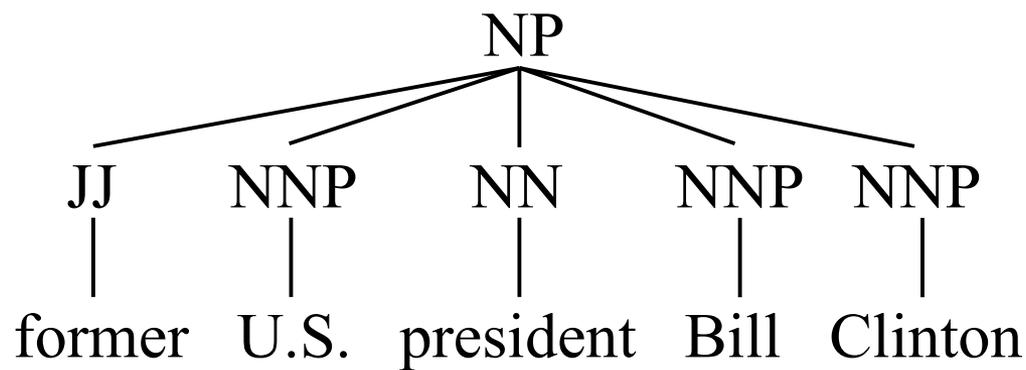
- Larger grammars: score difference erased

# Conclusions

- Very large linguistically motivated rule sets
  - No violating constituent bounds (Stat-XFER)
  - Multiple node alignments
  - Multiple decompositions (Hiero, GHKM)
  - Virtual nodes (< SAMT)
- More phrase pairs help improve scores
- Grammar filtering also matters

# Future Work

- Filtering to limit derivational ambiguity
- Filtering based on content of virtual nodes



- Reducing the size of the label set
  - Original: 1,577
  - With virtual nodes: 73,000

# References

- Chiang (2005), "A hierarchical phrase-based model for statistical machine translation," ACL
- Chiang (2010), "Learning to translate with source and target syntax," ACL
- Galley, Hopkins, Knight, and Marcu (2004), "What's in a translation rule?," NAACL
- Lavie, Parlikar, and Ambati (2008), "Syntax-driven learning of sub-sentential translation equivalents and translation rules from parsed parallel corpora," SSST-2
- Zollmann and Venugopal (2006), "Syntax augmented machine translation via chart parsing," WMT