# A Tookit for Visualizing the Coherence of Tree-based Reordering with Word Alignments

## Authors:

**Gideon Maillette de Buy Wenniger,**

**Maxim Khalilov,  Khalil Sima'an**

**Statistical Language Processing and Learning Lab at the Institute for Logic Language and Computation (ILLC), University of Amsterdam, the Netherlands**
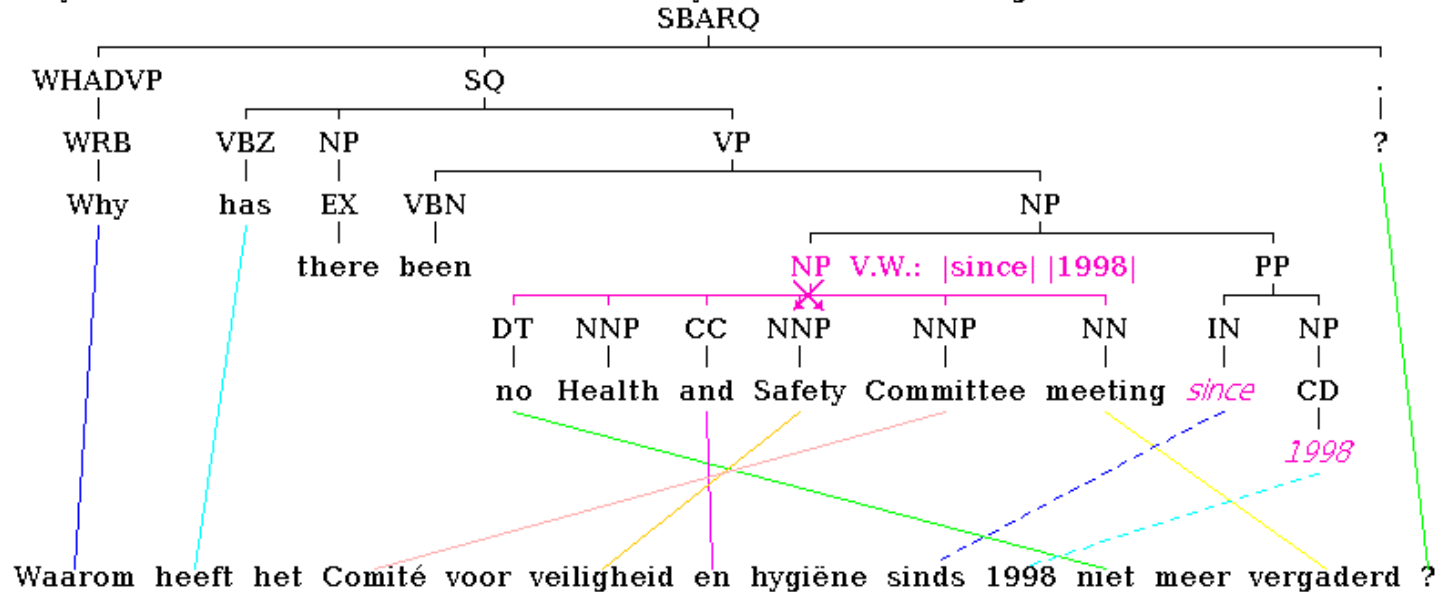
- Improving SMT:
  - Better models
  - Better learning methods and decoding

- Better  models:
  - More sensible alignments
  - Explicit language-specific reordering models
  - Adding all sorts of extra information

- This work: support search better models

- Data visualization facilitates SMT research

- Tree Structures: basis **syntactic SMT**

- *Un-cohesiveness* resulting from negation



- Is this big NP subtree appropriate?
  - For translation?
  - For reordering?
- Insight in **coherence** trees and alignments

# Basic Alignment Visualization

■ **Alignment between source and target sentence**
  ◆ **General: m-to-n mappings between words**
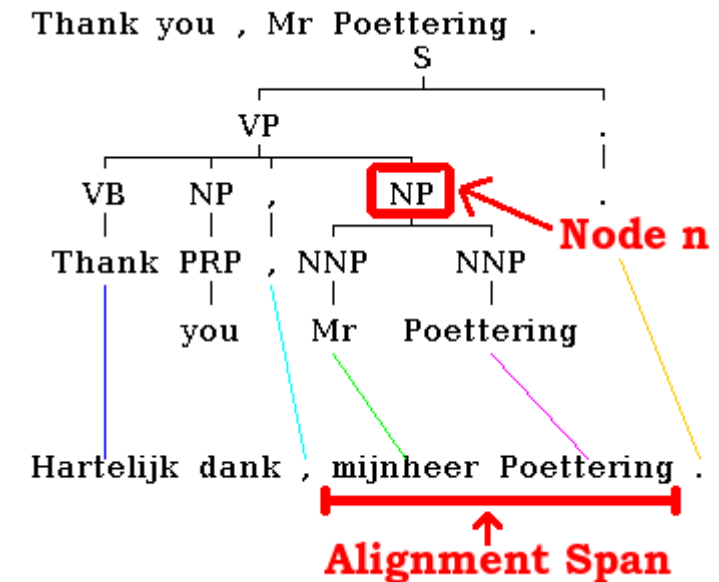
- Alignment mapping function

$$A(n) \rightarrow \{1, \cdots, m\}^*$$

- Span of target positions covered by subtree rooted at node **n**



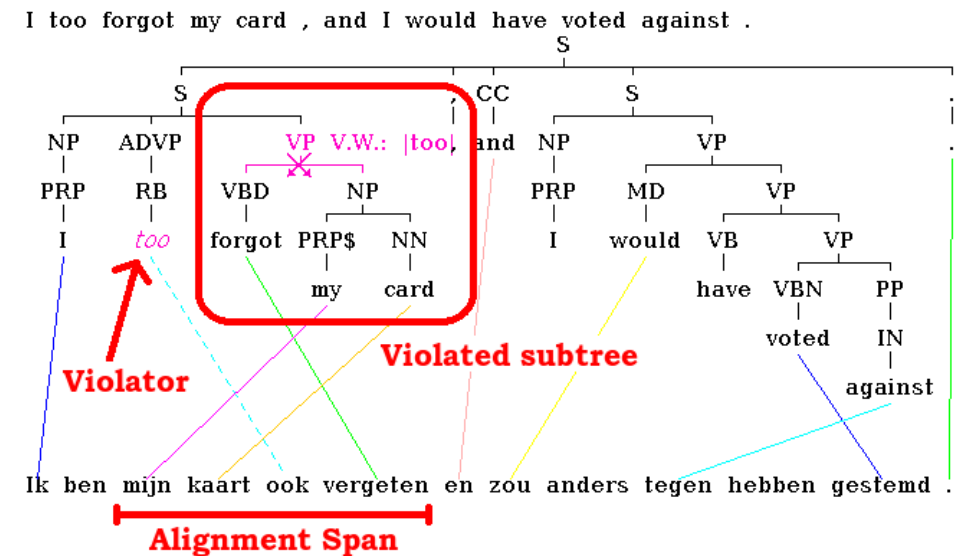**Definition 2.1** (Alignment Span)

$$\text{AlignmentSpan}(n) :=$$

$$[a_{n_{min}}, a_{n_{max}}] = \left[ \min_{x \in \text{LeafNodes}(n)} \left( \min_{a_{x'} \in A(x)} a_{x'} \right), \right.$$

$$\left. \max_{y \in \text{LeafNodes}(n)} \left( \max_{a_{y'} \in A(y)} a_{y'} \right) \right]$$

- Alignment  *Cohesive* nodes: source side syntactic phrase pair

- Un-cohesive / **Alignment Violation**: Two distinct subtrees align within same target range



**Definition 2.2** (Alignment Violation)

$$violates(n', n) := terminal(n') \wedge n' \notin descendants(n) \wedge$$
$$(AlignmentSpan(n) = [a_{n_{min}}, a_{n_{max}}]) \wedge$$
$$(a_{n_{min}} \leq A(n') \leq a_{n_{max}})$$

- **Ask** and **once** violate the alignment span of the S subtree

- Origin:
  *Inversion Transduction Grammars* (Wu, 1997):
  Bilingual Parsing

- Application   for Reordering:
  - Basic: child nodes binary tree may be inverted
  - General tree: permute child nodes arbitrarily
  - Restriction to constituency parse

- ***Monotonization:*** Reorder source to match target order

- **Formal definition precedence used**

**Definition 2.3** (Alignment Span Precedence)

$$\text{AlignmentSpan}(c1) = [a1_{\min}, a1_{\max}] <$$
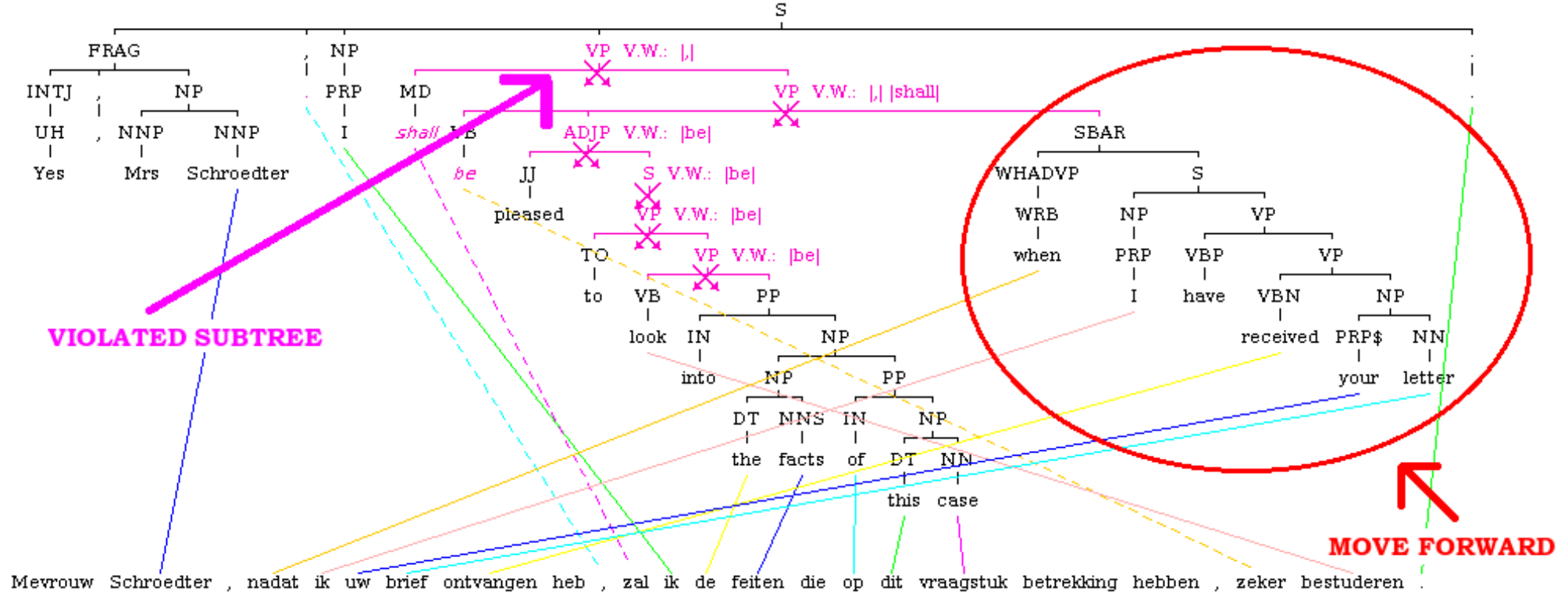$$\text{AlignmentSpan}(c2) = [a2_{\min}, a2_{\max}]$$
$$:= (a1_{\min} < a2_{\min}) \wedge (a1_{\max} < a2_{\min})$$

- **All positions covered by *c1* strictly precede those covered by *c2***

- ***Un-cohesiveness*** causes order problems

- **Right subtree moved forward, violated subtree not touched**

- Browsing through aligned sentences

- Insight into alignment mapping sub-trees

- Assess quality reordering tree-constrained ITG

- Get ideas for new tree-transduction operations

■ Visualization  tools
  ◆ *Cairo* (Smith and Jahr, 2000)
  ◆ *Yawat* (Germann, 2008)
  ◆ *Stockholm Tree Aligner (STA)* (Volk et al., 2007)

- Visualization alignment matrix

- Support manual annotation

- Dynamic highlighting

- ## Visualization parallel treebanks

- ## Focus on hand-annotated trees

- Toolkit targeted especially SMT people

- Focus on automatically generated resources and syntactic SMT

- Offers new functionality

- Goal: support SMT research

- Heuristics in subtree reordering

- Tree modification

- Alignment refinement

Burkett, D., J. Blitzer, , and D. Klein. Joint parsing and alignment with weakly synchronized grammars. In *Proc. of the 11th Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*, 2010.

Collins, M., P. Koehn, and I. Kučerová. Clause restructuring for statistical machine translation. In *Proc. of the 43rd Annual Meeting on Association for Computational Linguistics (ACL)*, pages 531–540, 2005.

Costa-jussà, M. R. and J. A. R. Fonollosa. Statistical machine reordering. In *Proc. of the 2006 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 70–76, 2006.

Galley, M., M. Hopkins, K. Knight, and D. Marcu. What's in a translation rule? In *Proc. of the Human Language Technology Conference and the North American Association for Computational Linguistics (HLT-NAACL)*, pages 273–280, 2004.

Germann, U. Yawat: yet another word alignment tool. In *Proc. of the 46th Annual Meeting of the Association for Computational Linguistics on Human Language Technologies (ACL-HLT)*, pages 20–23, 2008.

Khalilov, M. and K. Sima'an. A discriminative syntactic model for source permutation via tree transduction. In *Proc. of The Fourth Workshop on Syntax and Structure in Statistical Translation (SSST-4) at the 23rd International Conference on Computational Linguistics (COLING)*, pages – to appear, 2010.

Koehn, P., F.J. Och, and D. Marcu. Statistical phrase-based translation. In *Proc. of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology (HLT-NAACL)*, 2003.

Mariño, J. B., R. E. Banchs, J. M. Crego, A. de Gispert, P. Lambert, J. A. R. Fonollosa, and M. R.Costa-jussà. N-gram based machine translation. *Computational Linguistics*, 32(4):527–549, 2006.

Och, F. and H. Ney. The alignment template approach to statistical machine translation. *Computational Linguistics*, 30:417–449, 2004.

Pauls, A., D. Klein, D. Chiang, and K. Knight. Unsupervised syntactic alignment with inversion transduction grammars. In *Proc. of 11th Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL 2010)*, 2010.

Smith, Noah A. and Michael E. Jahr. Cairo: An alignment visualization tool. In *Proc. of the 2nd Conference on Language Resources and Evaluation (LREC)*, page 549–551, 2000.

Volk, M., J. Lundborg, and M. Mettler. Alignment tools for parallel treebanks. In *In Proc. of The Linguistic Annotation Workshop at the Association for Computational Linguistics (LAW-ACL)*, 2007.

Wang, W., J. May, K. Knight, and D. Marcu. Re-structuring, re-labeling and re-aligning for syntax-based machine translation. *Computational Linguistics*, 36:247–277, 2010.

Wu, D. Stochastic inversion transduction grammars and bilingual parsing of parallel corpora. *Computational Linguistics*, 23:377–404, 1997.

Xia, F. and M. McCord. Improving a statistical MT system with automatically learned rewrite patterns. In *Proc. of the 20th international conference on Computational Linguistics (COLING)*, pages 508–514, 2004.

Yamamoto, H., H. Okuma, and E. Sumita. Imposing constraints from the source tree on itg constraints for smt. In *Proc. of the Second Workshop on Syntax and Structure in Statistical Translation (SSST '08)*, page 1–9, 2008.