



Multiword Expression Identification with Recurring Tree Fragments and Association Measures

Proceedings of the 11th Workshop on Multiword Expressions

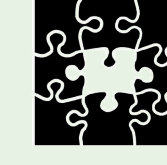
WG3

Parseme 5th GM
Iași, 23-24 September 2015

Federico Sangati
federico.sangati@gmail.com



Andreas van Cranenburgh
andreas.van.cranenburgh@huygens.knaw.nl



INSTITUTE FOR LOGIC,
LANGUAGE AND COMPUTATION



GOAL

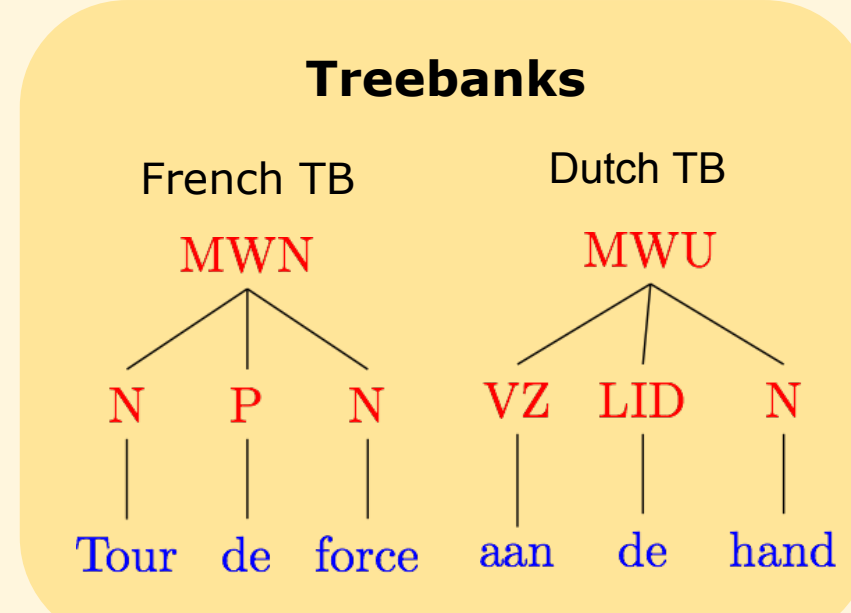
Investigate ways of **automatically discovering MWEs** in treebanks by searching for **recurring patterns**.

Trebank	Trees	Total Frags	Selected Frags
French	13K	274K	86K
Dutch	52K	536K	193K
English	500K	4.3M	2.8M

SUPERVISED

Input: Small MANUAL TREEBANKS with ANNOTATED MWE

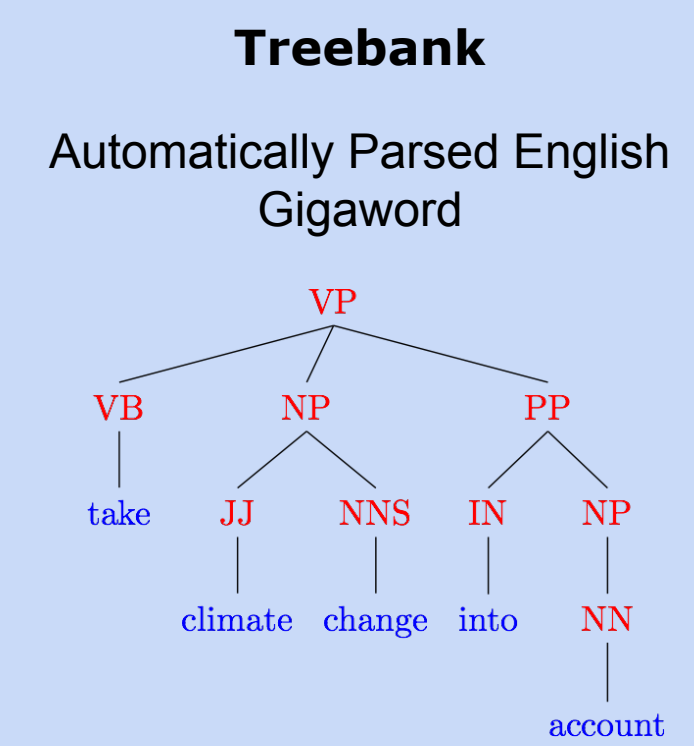
- Manually annotated MWEs.
- Only **contiguous** and **flat** MWEs



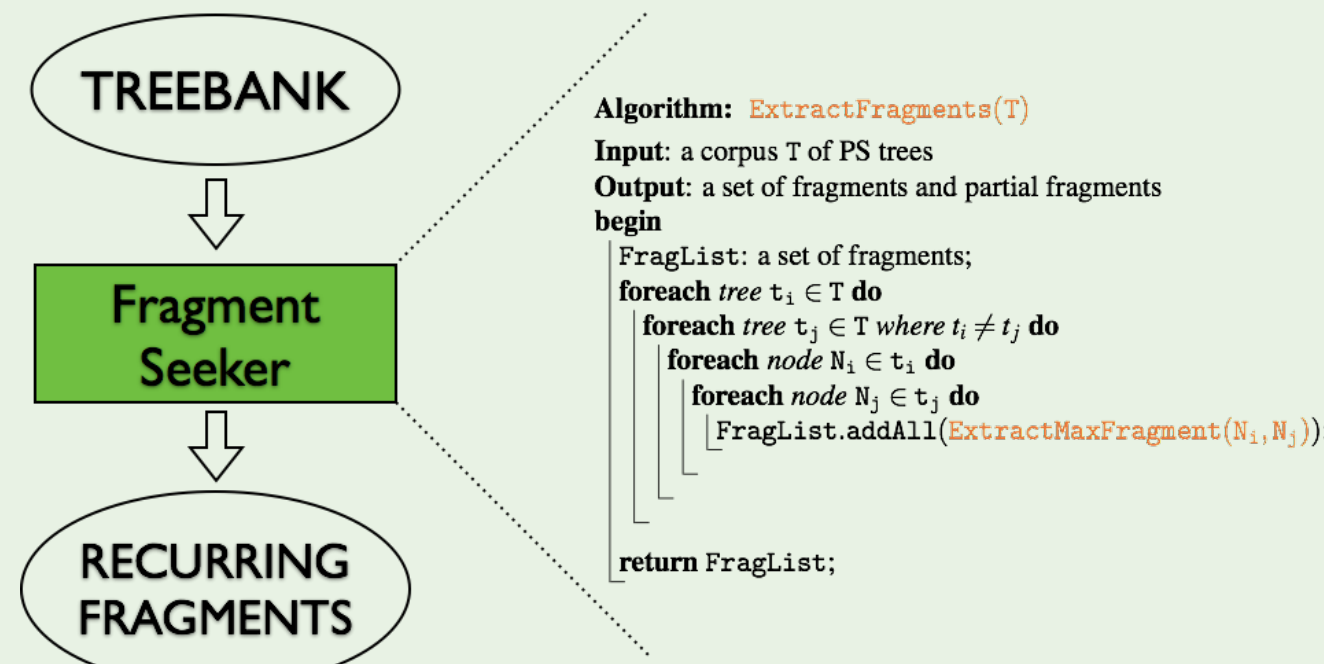
UNSUPERVISED

LARGE AUTOMATIC TREEBANK **without** ANNOTATED MWE

- Also **discontiguous** and **hierarchical** MWEs
- Not** manually annotate MWEs



RECURRING FRAGMENTS SEEKER



Parsing

Using **Disco-DOP: Tree Substitution Grammar**
(Implicit identification of MWEs as in Green et al. 2011)

Parser	F1	EX	MWE-F1
FRENCH			
Green et al. (2013): DP-TSG	76.9	16.0	71.3
Green et al. (2013): Stanford	79.0	17.6	70.5
disco-dop, 2DOP	79.3	19.9	71.9
DUTCH			
disco-dop, PCFG baseline	63.9	21.8	50.4
disco-dop, 2DOP	77.0	35.2	75.3

#gold	DP-TSG	Stanford	This work
MWN	457	65.7	64.8
MWADV	220	77.2	75.0
MWP	162	79.5	81.2
MWC	47	85.8	86.3
MWV	26	56.2	57.1
MWPRO	17	75.3	72.2
MWD	15	65.1	68.4
MWA	8	36.0	26.1
Total	955	71.3	70.5

RESULTS

Computing single-figure **F1 evaluation** of the 3 metrics, obtained by aggregating the top 1/5 candidates of each bin. For this evaluation, recall and precision are computed, with the gold set consisting of all the extracted lexicalized fragments with MWE gold tags.

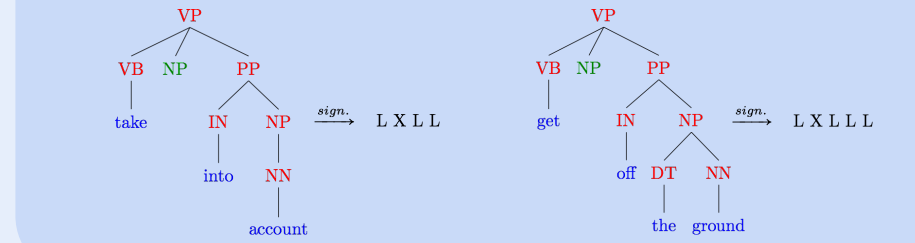
Trebank	PMI	LLR	SCI
French	33.0	32.3	45.8
Dutch	49.4	46.6	50.5

This **evaluation** is **not ideal**, as our method aims to **go beyond the small, contiguous MWE strings** annotated in the treebanks.

Manual inspection of the selected candidates reveals that many of them are MWEs, while not part of the gold standard.

MWEs Discovery

1. Partitioning Fragments in **signature** bins:



2. Ranking fragments based on

a. Association Measures

Log-Likelihood Ratio

$$LLR(S_1, \dots, S_n) = \log \frac{p(S_1, \dots, S_n)}{\sum_{\sigma \in CSP(S_1, \dots, S_n)} \prod_{\sigma} p(\sigma)}$$

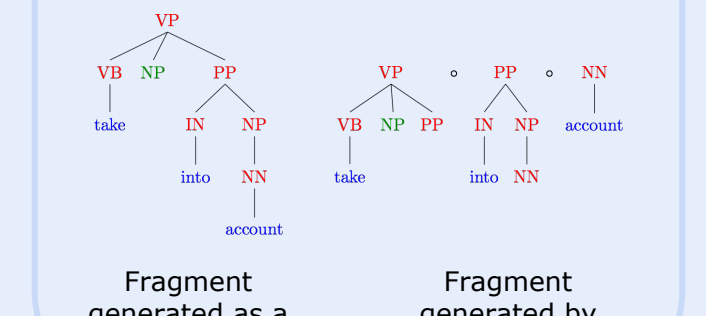
Pointwise Mutual Information

$$PMI(S_1, S_2, \dots, S_n) = \log \frac{p(S_1, S_2, \dots, S_n)}{\prod_{i=1}^n p(S_i)}$$

b. Syntax Compositionality Index

SCI(frag) = $\log \frac{p(\text{frag})}{\text{inside}(\text{frag})}$

High when a fragment is often seen as a single block, low when it is typically generated by smaller units



Fragment generated as a single block

Fragment generated by smaller units