

Automatic meta-grammar extraction

Matías Guzmán Naranjo

21.02.2018, Düsseldorf

Two steps

- \wedge : split and cut trees into fragments
- \vee : rearrange fragments into meaningful groups

\wedge

Splitting trees

What does not seem to work: combinatoric approaches.

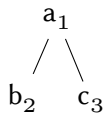
- treat trees like sets.
- find the cover set of each tree (all possible decompositions)
- try to find a reasonable set of decompositions in the chaos of trillions of possible decompositions.

Splitting trees

What does seem to work: use small trees to cut bigger trees.

- start with a (possibly empty) set of tree fragments
- sort all trees by size (number of edges, small first)
- starting with the smallest tree, try to cut it into pieces found in the set of tree fragments (sorted from largest to smallest)
- add the 'left overs' (possibly the whole tree) to the inventory of tree fragments
- **(add the tree itself to the inventory of tree fragments)**
- repeat until you've cut all your trees

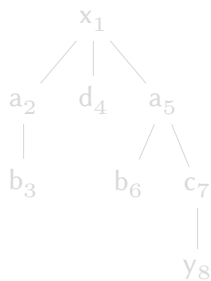
Cutting trees



—



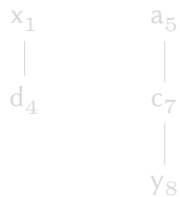
⇒



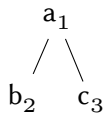
—



⇒



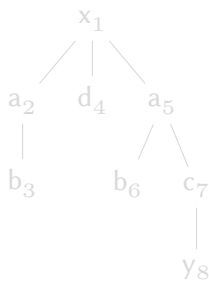
Cutting trees



—



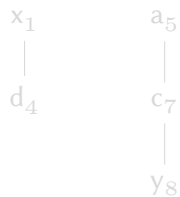
⇒



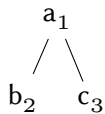
—



⇒



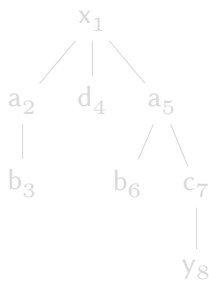
Cutting trees



—



⇒



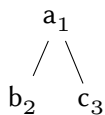
—



⇒



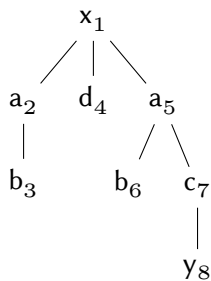
Cutting trees



—



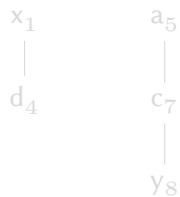
⇒



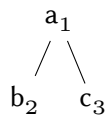
—



⇒



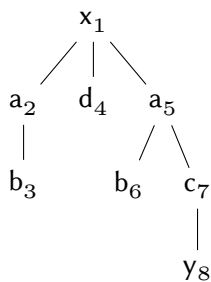
Cutting trees



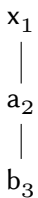
—



⇒



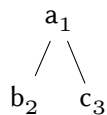
—



⇒



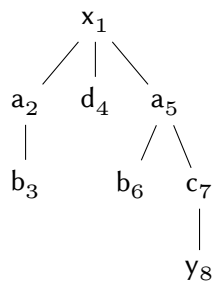
Cutting trees



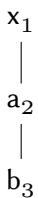
—



⇒



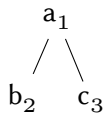
—



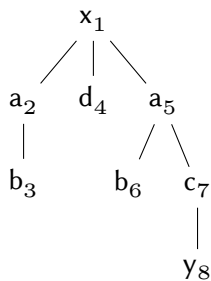
⇒



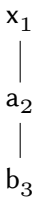
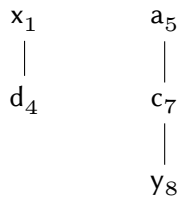
Cutting trees



—

 \Rightarrow 

—

 \Rightarrow 

Shortcomings

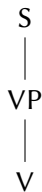
This approach has some clear shortcomings

- The decompositions 'lose' linear precedence
- There is no good way of naming the resulting fragments
- No abstraction across trees

Seeds

An interesting aspect of this approach is that it allows us to start the system with user defined seeds.

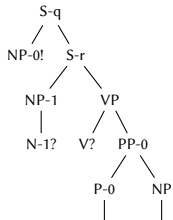
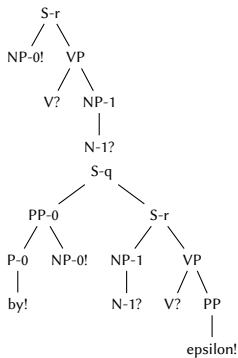
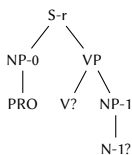
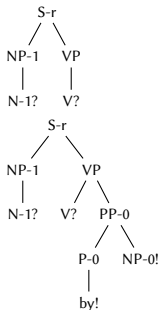
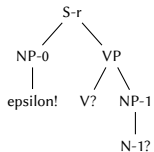
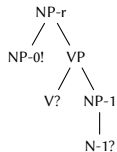
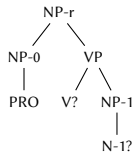
We can include things like the vp-spine as a basic fragment:



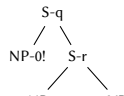
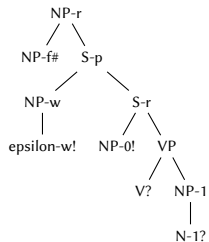
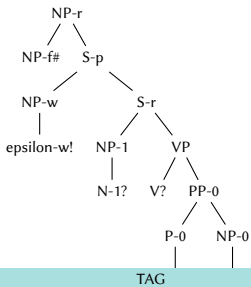
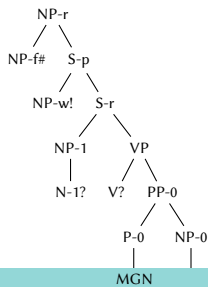
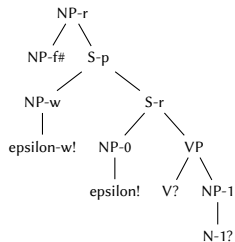
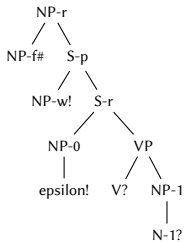
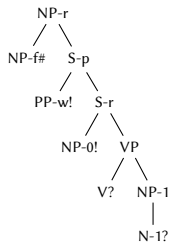
This guarantees that all trees where the vp-spine is present, will be decomposed as having this fragment.

The final factorization has 719 unique fragments (for 1300 trees), but many of these are reused trees.

Tnx0VN1



Tnx0VN1

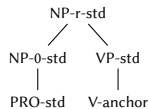


Tnx0VN1

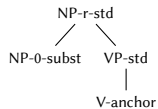
alphaGnx0VN1-PRO:	alphaGnx0V-PRO , VP-std->NP-1-std NP-1-std->N-1-anchor
alphaGnx0VN1:	alphaGnx0V , VP-std->NP-1-std NP-1-std->N-1-anchor
alphaInx0VN1:	alphaInx0V , VP-std->NP-1-std NP-1-std->N-1-anchor
alphaN1V:	vp-spine , S-r-std->NP-1-std NP-1-std->N-1-anchor
alphaN1Vbynx0:	alphaN1V , VP-std->PP-0-std PP-0-std->P-0-std PP-0-std->NP-0-subst P-0-std->by-flex6»7
alphanx0VN1-PRO:	alphanx0V-PRO , VP-std->NP-1-std NP-1-std->N-1-anchor
alphanx0VN1:	alphanx0V , VP-std->NP-1-std NP-1-std->N-1-anchor
alphapW0N1Vbynx0:	S-q-std->PP-0-std S-q-std->S-r-std PP-0-std->P-0-std PP-0-std->NP-0-subst P-0-std->by-flex2»3-4»5, alphaN1V , VP-std->PP-std PP-std->epsilon-flex
alphaW0N1Vbynx0:	VP-std->PP-0-std PP-0-std->P-0-std PP-0-std->NP-std P-0-std->by-flex NP-std->epsilon-flex8»9, alphaN1V , S-q-std->NP-0-subst S-q-std->S-r-std2»3
alphaW0nx0VN1:	alphaW0nx0V , VP-std->NP-1-std NP-1-std->N-1-anchor
betaN0N1Vbynx0:	VP-std->PP-0-std PP-0-std->P-0-std PP-0-std->NP-0-std P-0-std->by-flex NP-0-std->epsilon-flex10»11, alphaN1V , NP-r-std->NP-f-foot NP-r-std->S-p-std S-p-std->NP-w-subst S-p-std->S-r-std2»3-4»5
betaN0nx0VN1:	betaN0nx0V , VP-std->NP-1-std NP-1-std->N-1-anchor
betaNc0N1Vbynx0:	NP-r-std->NP-f-foot NP-r-std->S-p-std S-p-std->NP-w-std S-p-std->S-r-std NP-w-std->epsilon-w-flex2»3-4»5, VP-std->PP-0-std PP-0-std->P-0-std PP-0-std->NP-0-std P-0-std->by-flex NP-0-std->epsilon-flex10»11, alphaN1V
betaNc0nx0VN1:	betaNc0nx0V , VP-std->NP-1-std NP-1-std->N-1-anchor
betaNcnx0VN1:	betaNcnx0V , VP-std->NP-1-std NP-1-std->N-1-anchor
betaNpxnx0VN1:	betaNpxnx0V , VP-std->NP-1-std NP-1-std->N-1-anchor

Tnx0VN1 fragments

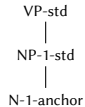
alphaGnx0V-PRO



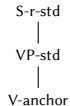
alphaGnx0V



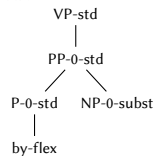
NN



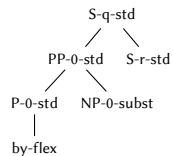
vp-spine

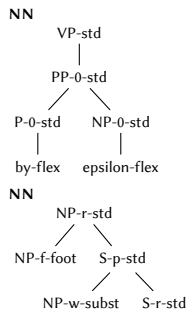
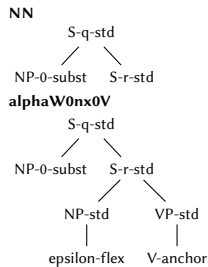
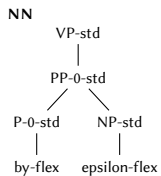
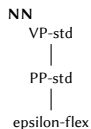


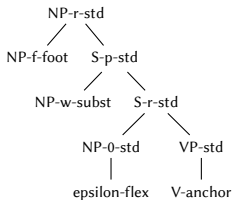
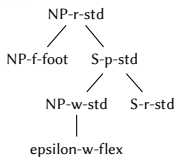
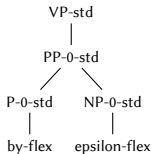
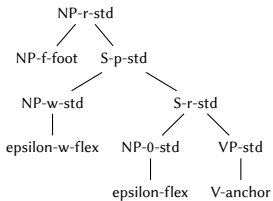
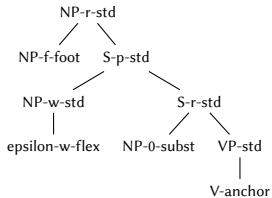
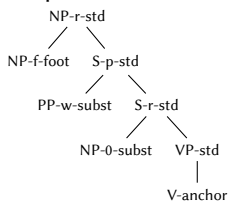
NN



NN



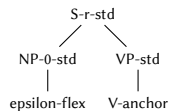


betaN0nx0V**NN****NN****betaNc0nx0V****betaNcnx0V****betaNpxnx0V**

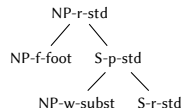
The bigger pices are themselves factorized

So, just to show you, in the previous slide the fragment **betaN0nx0V** gets factorized as:

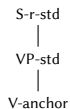
alphaInx0V



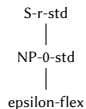
NN



Vspine



NN



Recompiling

When this factorization is exported back to xmg we get back almost all of the original trees (there are still some mistakes).

however:

- We have to explicitly define each tree
- add all missing node ordering constraints which are lost in the decomposition
- There are no \vee relations (!)

\vee

How to proceed?

Suppose the following tree decompositions:

- $t_1 = (a, b, c, d)$
- $t_2 = (a, b, c, e)$

Trivially, we can see that $a \wedge b \wedge c \wedge (d \vee e)$ produces both trees.

But if we have:

- $t_1 = (a, b, c, d)$
- $t_2 = (a, b, c, e)$
- $t_3 = (a, f, c, e)$
- $t_4 = (a, f, c, d)$
- $t_5 = (g, b, c, e)$

What to do?!

In the intransitive family

alphanx0V-PRO	⇒	(ae, ag)
betaNc0nx0V	⇒	(al , aj)
betaVintransn	⇒	(ao)
alphaGnx0V-PRO	⇒	(ac)
alphaDnx0V	⇒	(ab)
alphaW0nx0V	⇒	(ae, ai)
alphaInx0V	⇒	(ae, af)
alphanx0V	⇒	(ae, ah)
betaNcnx0V	⇒	(al , am)
betaN0nx0V	⇒	(aj, ak)
betaNpxnx0V	⇒	(an, am)
alphaGnx0V	⇒	(ad)

In the transitive family

$\beta N 0 n x 1 V b y n x 0 \Rightarrow a z, a p, b c$
 $\beta N 0 n x 0 V n x 1 \Rightarrow a y, a g$
 $\alpha p h a n x 1 V - P R O \Rightarrow a o$
 $\alpha p h a W 0 n x 0 V n x 1 \Rightarrow a s, a g$
 $\alpha p h a W 0 n x 1 V b y n x 0 \Rightarrow a p, a t, a u$
 $\alpha p h a p W 0 n x 1 V b y n x 0 \Rightarrow a p, a q, a r$
 $\alpha p h a n x 1 V b y n x 0 - P R O \Rightarrow a o, a k$
 $\beta N c 0 n x 0 V n x 1 \Rightarrow b h, a g$
 $\beta N c n x 1 V b y n x 0 \Rightarrow b i, a k$
 $\alpha p h a n x 1 V b y n x 0 \Rightarrow a k, a p$
 $\alpha p h a I n x 0 V n x 1 \Rightarrow a l, a g$
 $\alpha p h a n x 0 V n x 1 \Rightarrow a n, a g$
 $\beta N 1 n x 1 V b y n x 0 \Rightarrow b e, a k$
 $\beta N p x n x 0 V n x 1 \Rightarrow b l, a g$
 $\beta N c 1 n x 0 V n x 1 \Rightarrow b j, b d$
 $\beta N p x n x 1 V \Rightarrow b m$
 $\alpha p h a W 1 n x 1 V b y n x 0 \Rightarrow a x, a k$
 $\beta N c 1 n x 1 V b y n x 0 \Rightarrow b k, a k$
 $\beta N 1 n x 1 V \Rightarrow b e$

$\beta N b y n x 0 n x 1 V b y n x 0 \Rightarrow a p, b f, b g$
 $\alpha p h a n x 0 V n x 1 - P R O \Rightarrow a m, a g$
 $\alpha p h a A V \Rightarrow a b$
 $\alpha p h a W 1 n x 1 V \Rightarrow a x$
 $\alpha p h a W 1 n x 0 V n x 1 \Rightarrow a n, a v, a w$
 $\alpha p h a G n x 1 V b y n x 0 - P R O \Rightarrow a i, a k$
 $\beta N 1 n x 0 V n x 1 \Rightarrow a z, a n, b d$
 $\alpha p h a G n x 0 V n x 1 \Rightarrow a h, a g$
 $\alpha p h a G n x 1 V \Rightarrow a j$
 $\alpha p h a G n x 1 V - P R O \Rightarrow a i$
 $\beta N c 1 n x 1 V \Rightarrow b k$
 $\alpha p h a D n x 0 V n x 1 \Rightarrow a c, a d, a e$
 $\beta N c n x 1 V \Rightarrow b i$
 $\alpha p h a G n x 1 V b y n x 0 \Rightarrow a j, a k$
 $\alpha p h a G n x 0 V n x 1 - P R O \Rightarrow a f, a g$
 $\beta N c 1 n x 1 V \Rightarrow b n$
 $\beta N c 0 n x 1 V b y n x 0 \Rightarrow b i, b c$
 $\beta N c n x 0 V n x 1 \Rightarrow b j, a g$
 $\alpha p h a n x 1 V \Rightarrow a p$
 $\beta N p x n x 1 V b y n x 0 \Rightarrow b m, a k$

Fragment frequency

ag	⇒	10
ak	⇒	9
ap	⇒	6
bi	⇒	3
an	⇒	3
bd	⇒	2
aj	⇒	2
be	⇒	2
ao	⇒	2
bm	⇒	2
az	⇒	2
ax	⇒	2
ai	⇒	2
bc	⇒	2
bk	⇒	2
bj	⇒	2

A trivial ∨ factorization

$ak \wedge (ao \vee bi \vee ap \vee be \vee ax \vee bk \vee ai \vee aj \vee bm)$

$ap \wedge ((bc \wedge az) \vee (at \wedge au) \vee (ar \wedge aq) \vee (bg \wedge bf))$

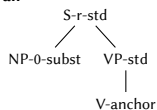
$ag \wedge (ay \vee as \vee bh \vee al \vee an \vee bl \vee am \vee ah \vee af \vee bj)$

$an \wedge ((av \wedge aw) \vee (bd \wedge az))$

$$\text{an} \wedge ((\text{av} \wedge \text{aw}) \vee (\text{bd} \wedge \text{az}))$$

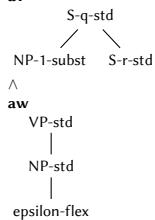
$$\Rightarrow \text{alphaW1nx0Vnx1}, \text{betaN1nx0Vnx1}$$

an



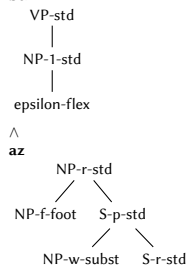
alphaW1nx0Vnx1

av



betaN1nx0Vnx1

bd



Many problems with this approach:

- We see that what both trees have in common is that they both have a subject
- Other trees with subjects are being ignored because they use larger trees as fragments
- Not very interesting linguistically (?)

Some alternatives (not tried so far)

- Try to infer 'families' with some distance based clustering
- (dis)Prefer more frequent trees as pivots for the \forall s
- Use features to only group together fragments which are identical.
- All of the above (?)

Factorizing feature structures

We take a similar approach (maybe?)

For a single fragment type (e.g. all cases of the vp-spine):

- Find the fragment instances with largest feature structures that are a subset of the feature structures of other fragment instances
- Decompose those fragment instance as being made up of two feature structures
- repeat until done

Problems

This has a series of problems, however:

- It is not clear that the decompositions make much sense (though maybe?)
- Decompositions seem to remain mostly binary
- Decompositions seem to mostly be ((ft1) (features))
- There is no natural way to sort the porcess, thus many different possible results.
- Maybe the feature decomposition should be linked is some other way to the tree splitting process?

That's it...