

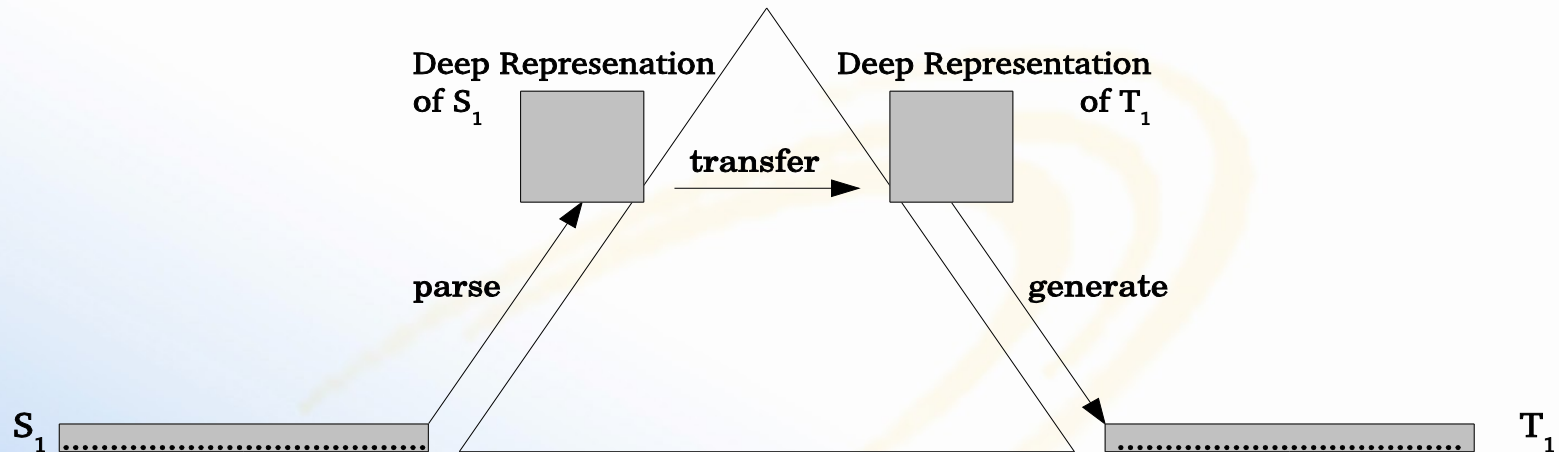


Probabilistic Transfer-based MT

Yvette Graham, Josef van Genabith
National Centre for Language Technology
School of Computing
Dublin City University

Transfer-based Machine Translation

National Centre for Language Technology

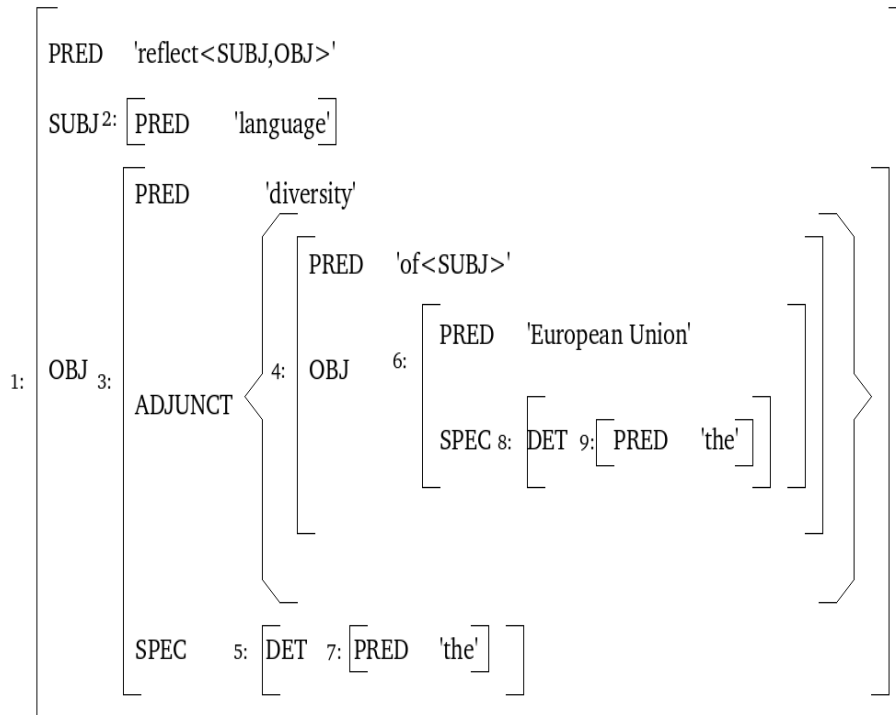


- Transfer-rule Induction Algorithm & Packed Rule Data Structure
- Probabilistic Decoder for Transfer-based MT

Why LFG F-structure?

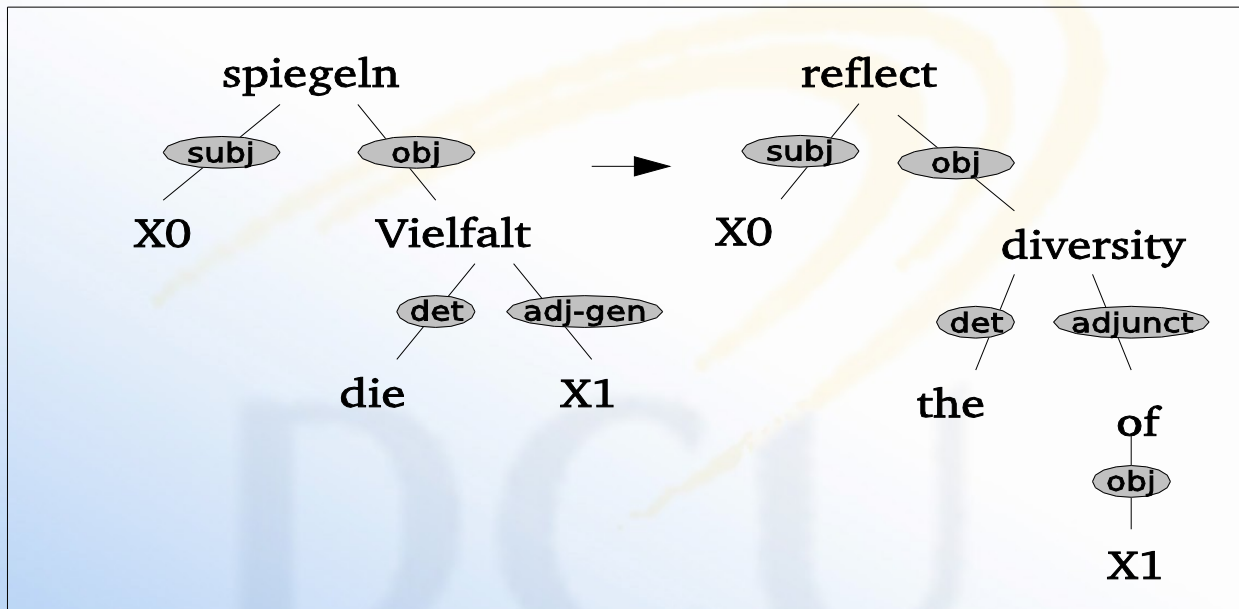


“Languages reflect the diversity of the European Union.”



- encode **dependencies** between the content words of a sentence
- encode the **linguistic information** necessary to generate the TL surface form sentence
- provide an abstract linguistic structure for forming **linguistically motivated generalizations** about how to translate well from one language to another
- abstracting away from surface form to f-structure produces **more general rules** that can be applied to multiple surface form examples – reducing data sparseness problems
- reliable and accurate linguistic **resources** for parsing and generation are vital to the overall success of a transfer-based system

Example F-structure Transfer Rules



X0 spiegeln(t) die Vielfalt **X1** wider. → **X0** reflect(s) the diversity of **X1**.

Transfer Rule Induction

National Centre for Language Technology



- Exponential number of possible transfer rules within a single f-structure pair
- Lots of those rules are complete rubbish!

The Plan: Constrain rule induction to eliminate unwanted rules

DCU

Constraining Rule Induction



- Contiguity Constraint (Riezler & Maxwell 06)
- Cross-structural Consistency Constraint:
 - Use a 1-1 set of alignments between nodes & allow any SL-TL node pair with equivalent non-empty sets of *aligned descendants* be a transfer-rule root.

Aligned Descendent of a SL local f-structure n:

Any descendent of n that is aligned with a TL f-structure.

If n itself is aligned with any TL local f-structure, then n is considered an aligned descendent of itself.

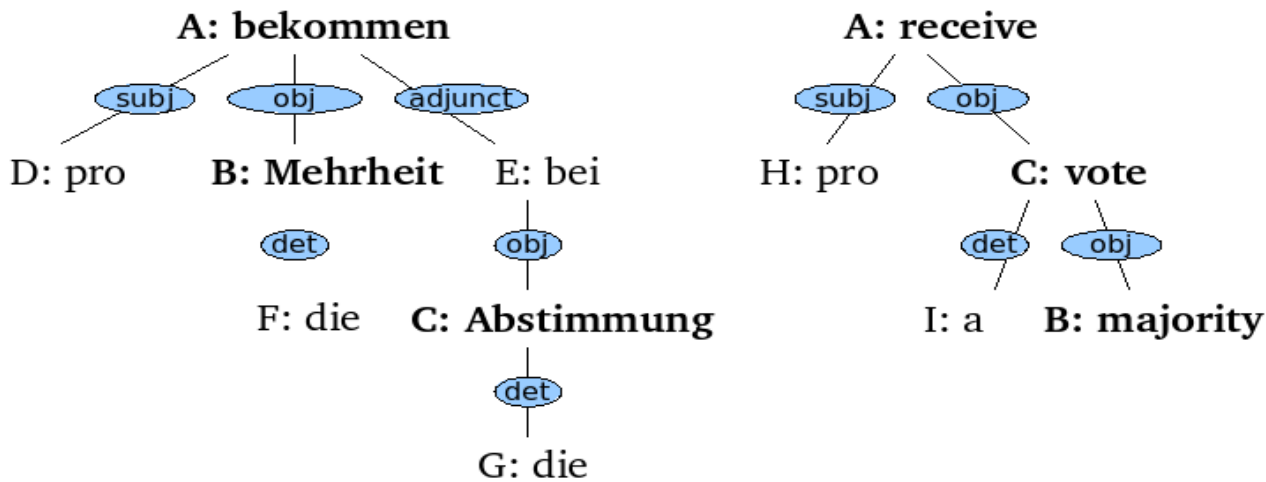
Note: by this definition the actual root nodes of every f-structure training pair is a rule root.

Example



Er hat die Mehrheit bei der Abstimmung bekommen.

It received a majority vote.






A: {A,B,C}
B: {B}
C: {C}
D: {}
E: {}
F: {}
G: {}

Rule Roots:
A,A
B,B

A: {A,B,C}
B: {B}
C: {B,C}
H: {}
I: {}

Induced Rules



<p> bekommen sub obj adjunct pro X0 bei obj Abstimmung det die </p> <p style="text-align: center;">↔</p> <p> receive sub obj pro vote det obj a X0 </p> <p>Er hat X bei der Abstimmung bekommen. ↔ It received a X vote.</p>	
<p> Mehrheit ↔ majority det die die Mehrheit ↔ majority </p>	
<p> Abstimmung ↔ vote det die a obj majority </p> <p> die Abstimmung ↔ a majority vote. </p>	

Putting Variables into Transfer Rules

National Centre for Language Technology



Variable Constraints:

1. The root of a transfer rule may never be a variable.
2. For any non-root node:
 - Iff it is aligned with a TL node it can be a variable

DCU



Example F-structure Transfer Rules

<p>a.</p> <p>spiegeln → reflect</p> <p>subj obj → subj obj</p> <p>X0 X1 X0 X1</p>	<p>b.</p> <p>spiegeln → reflect</p> <p>subj obj → subj obj</p> <p>Sprache X0 language X0</p>
<p>c.</p> <p>spiegeln → reflect</p> <p>subj obj → subj obj</p> <p>X0 Vielfalt X0 diversity</p> <p>det adj-gen det adjunct</p> <p>die X1 the X1</p>	<p>d.</p> <p>spiegeln → reflect</p> <p>subj obj → subj obj</p> <p>Sprache Vielfalt language diversity</p> <p>det adj-gen det adjunct</p> <p>die X1 the X1</p>
<p>e.</p> <p>spiegeln → reflect</p> <p>subj obj → subj obj</p> <p>X0 Vielfalt X0 diversity</p> <p>det adj-gen det adjunct obj</p> <p>die X1 the of X1</p>	<p>f.</p> <p>spiegeln → reflect</p> <p>subj obj → subj obj</p> <p>Sprache Vielfalt language diversity</p> <p>det adj-gen det adjunct obj</p> <p>die X1 the of X1</p>

Example F-structure Transfer Rules

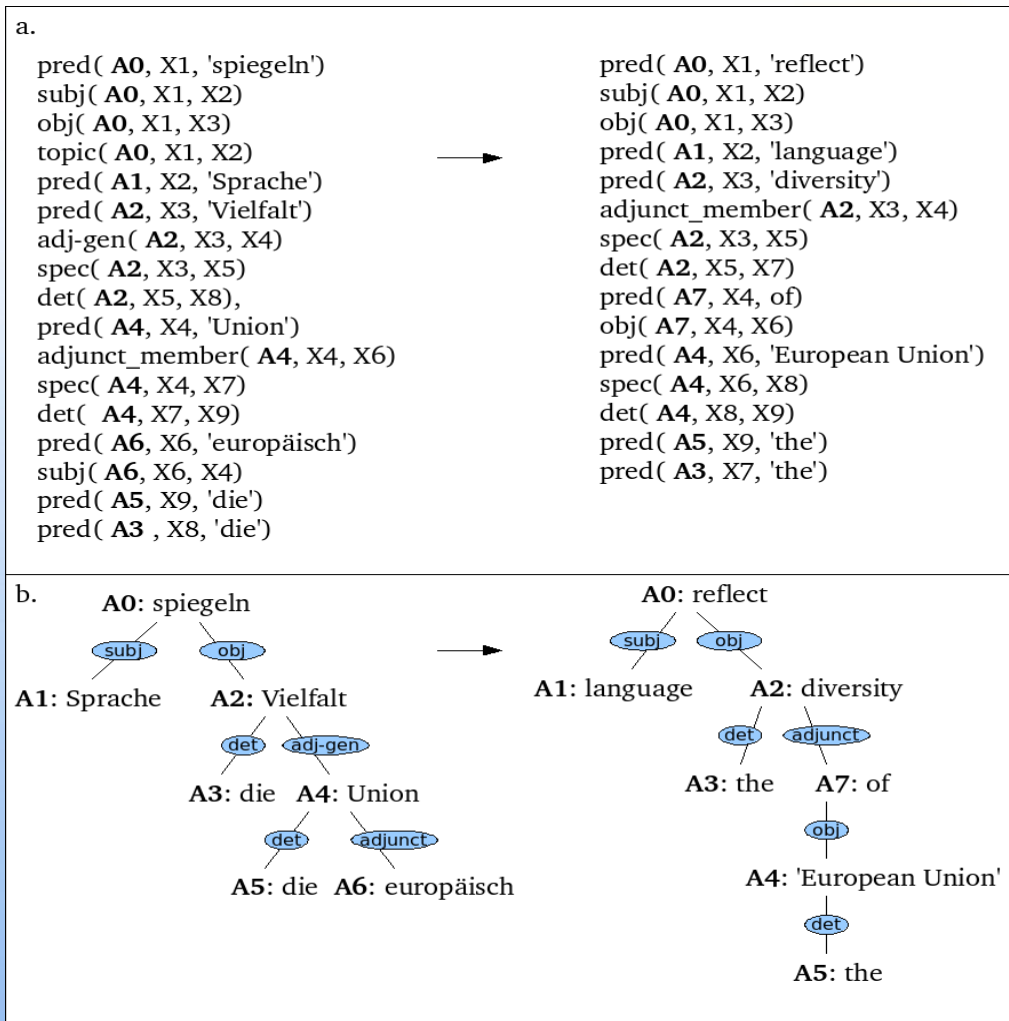


<p>g.</p>	<p>h.</p>
<p>i.</p> <p>Sprache → language</p>	<p>j.</p>
<p>k.</p>	<p>l.</p>
<p>m.</p>	<p>n.</p>

Packed Transfer Rule Data Structure



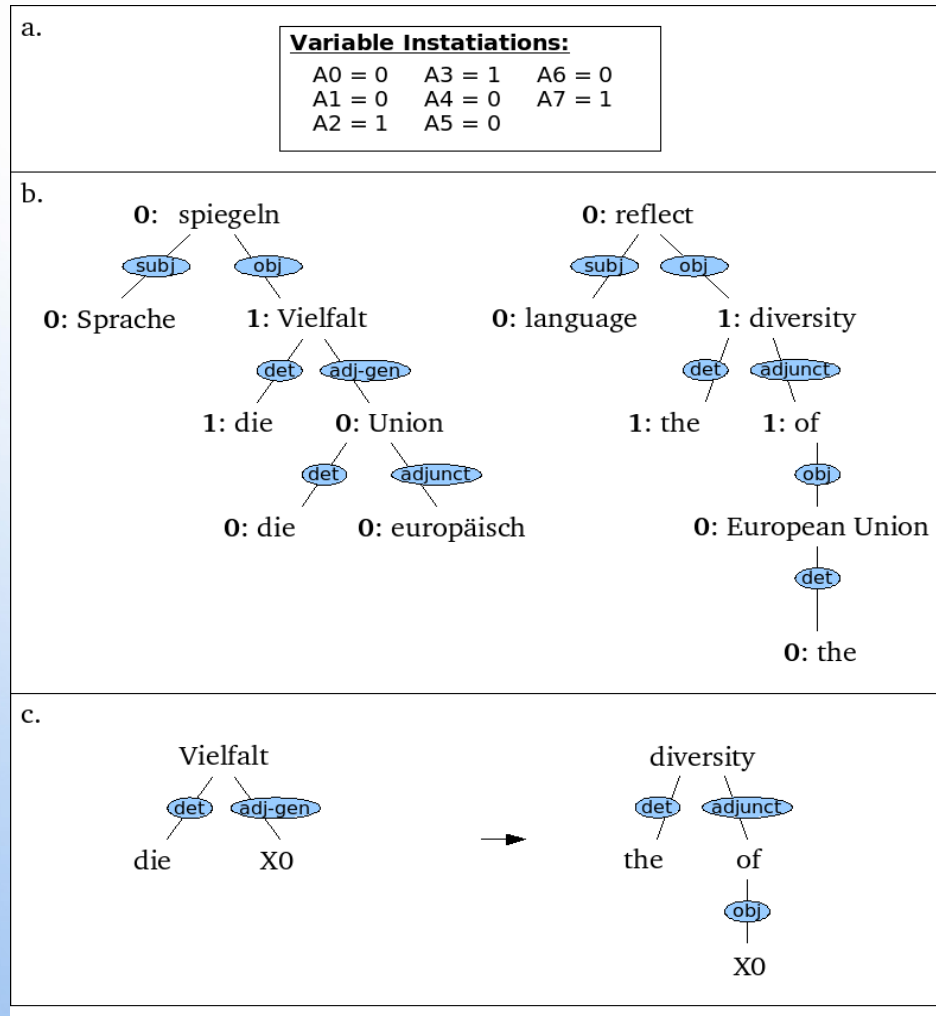
- Encode all possible rules that can be induced from a SL TL f-structure pair given these constraints in a single packed representation



Example Unpacking Rule



- Unpack a rule by assigning true or false to certain context variables



- To unpack all possible rules assign all possible combinations of true and false to the context variables

Experimental Evaluation of Packed Data Structure

National Centre for Language Technology



- Extracted all rules from 217,666 f-structure pairs of length 5-15 of German-English Europarl
- Unpacked 10 sets of rules – each set consisting of rules induced from randomly selected sets of 1000 training pairs, on average 23955 rules per set

Set	No. Rules	Disk Space		Write Time		Load Time		Retrieving/Unpacking	
		Enum	Packed	Enum	Packed	Enum	Packed	Enum	Packed
1	24121	96.37M	7.2M	144s	128s	211s	17s	2s	3s
2	24486	98.89M	7.16M	145s	127s	215s	19s	2s	3s
3	23650	93.58M	7.17M	142s	133s	200s	18s	2s	2s
4	23882	96.83M	7.22M	149s	118s	210s	18s	2s	2s
5	24146	98.03M	7.15M	148s	128s	212s	17s	2s	3s
6	23355	91.75M	7.1M	140s	128s	198s	21s	2s	3s
7	23620	94.55M	7.21M	141s	142s	204s	18s	2s	2s
8	23687	94.02M	7.11M	137s	124s	201s	17s	1s	3s
9	23534	94.95M	7.12M	142s	120s	204s	17s	1s	3s
10	25069	100.66M	7.26M	152s	231s	219s	19s	2s	2s
Average	23955	95.96M	7.17M	144s	137.9s	207.4s	18.1s	1.8s	2.6s
All Rules Estimate	5214189	20.4G	1.52G	8h43m	8h20m	12h32m	1h06m	6m31s	9m24s

Chart-based Decoder

National Centre for Language Technology



Input: a pre-compiled chart for the source language f-structure containing

Output: TL f-structure or n-best TL f-structures

- Top-down beam search of transfer chart
- Log-linear model to combine feature scores
- MERT to adjust weights (still to be integrated into decoder)

DCU

Future Work

National Centre for Language Technology



- Fully integrate **MERT** into decoder
- Implement a **faster search algorithm** for decoding
- **Factor atomic features** and values to investigate what level of rule specificity is best
- Automatically **learn where atomic features get their values** – SL v TL

DCU