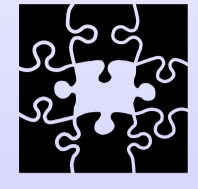


# A PROBABILISTIC GENERATIVE MODEL FOR AN INTERMEDIATE CONSTITUENCY-DEPENDENCY REPRESENTATION

Federico Sangati (f.sangati@uva.nl)  
ILLC, University of Amsterdam



INSTITUTE FOR LOGIC, LANGUAGE AND COMPUTATION



Netherlands Organisation for Scientific Research

**Abstract:** We present a probabilistic model extension to the Tesnière Dependency Structure (TDS) framework formulated in (Sangati and Mazza, 2009). This representation incorporates aspects from both constituency and dependency theory. In addition, it makes use of junction structures to handle coordination constructions. We test our model on parsing the English Penn WSJ treebank using a re-ranking framework. This technique allows us to efficiently test our model without needing a specialized parser, and to use the standard evaluation metric on the original Phrase Structure version of the treebank.

## The TDS representation: Tesnière Dependency Structures

**Words:**

- Content words (activities, encourage, ...)
- Functional words (the, that, for, ...)
- Punctuation words (, " " ...)

**Blocks:**

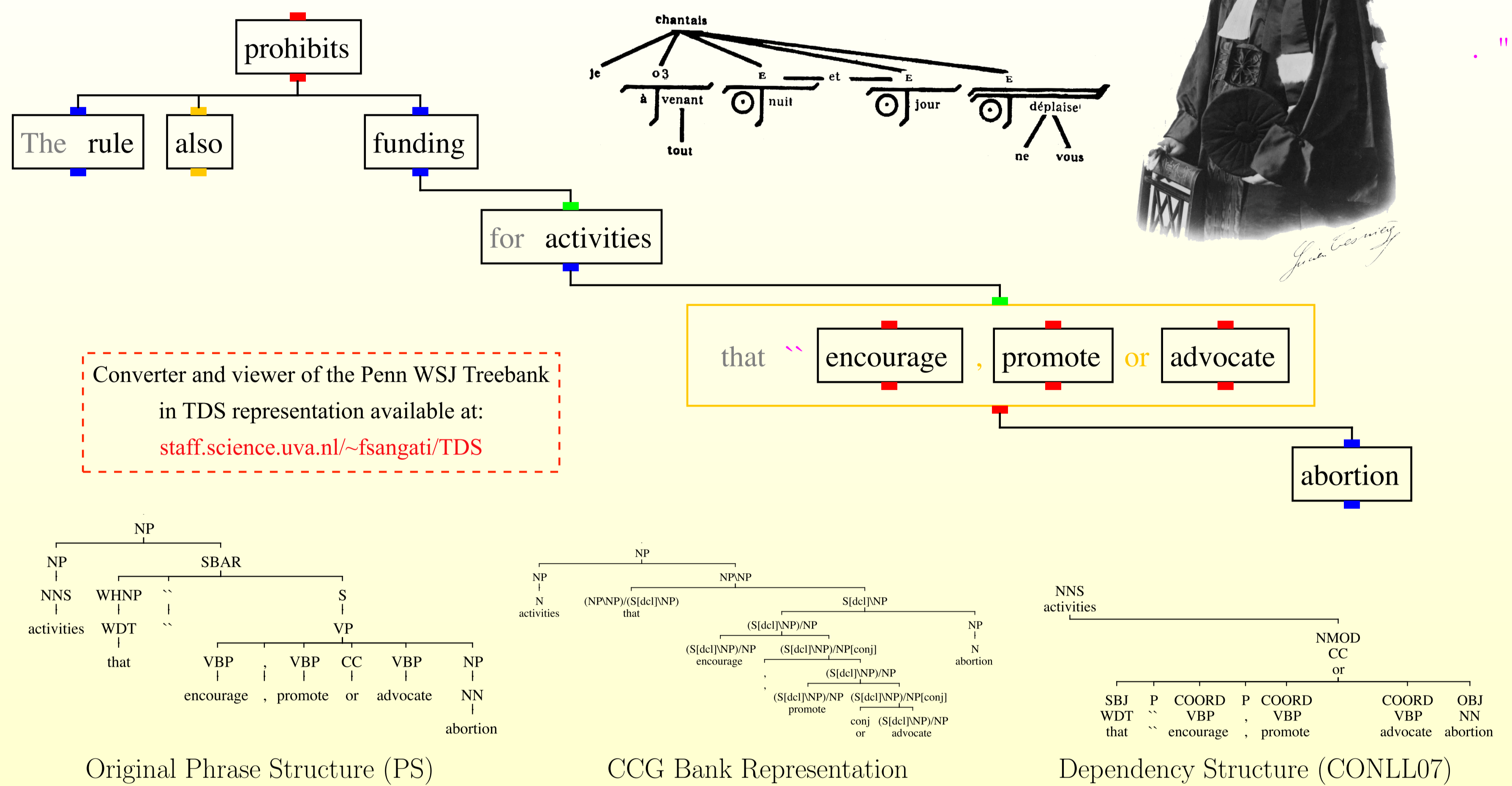
- Standard blocks
- Junction blocks

**Categories:**

- Verbs
- Adverbs
- Nouns
- Adjectives

Derived category: for activities

Original category: for activities

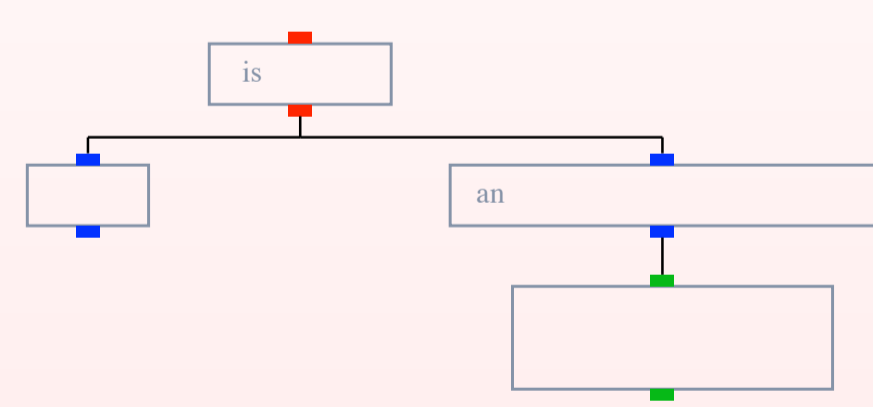


## A Probabilistic Generative Model for TDS

### The Model: 3 Generative Processes + PoS-Tagging & Chunking

#### 1. Block Generation

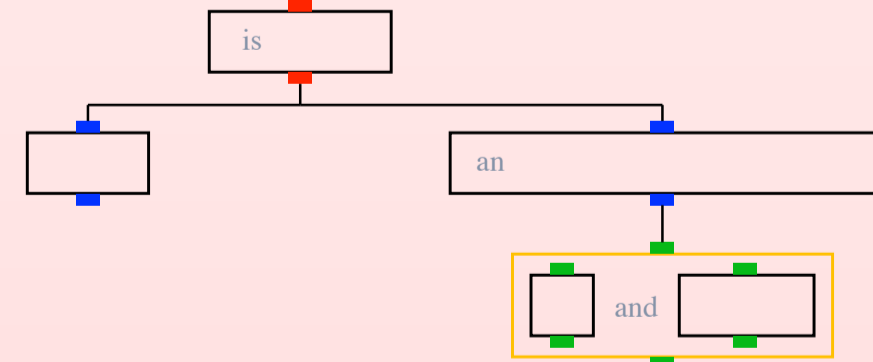
$$P_{BGM}(S) = \prod_{B \in \text{dependentBlocks}(S)} P(B|\text{parent}(B), \text{direction}(B), \text{leftSibling}(B))$$



B	Parent	Direction	Left Sibling
VER, VER, is	-	-	-
NOU, NOU	VER	Left	-
STOP	VER	Left	NOU, NOU
STOP	VER	Inner	-
NOU, NOU, an	VER	Right	-
STOP	VER	Right	NOU, NOU, an
STOP	NOU	Left	-
STOP	NOU	Right	-
STOP	NOU	Inner	-
STOP	NOU	Left	-
ADJ, ADJ	NOU	Inner	-
STOP	NOU	Inner	ADJ, ADJ
STOP	NOU	Right	-
3x STOP	ADJ	Left	-
3x STOP	ADJ	Inner	-
3x STOP	ADJ	Right	-

#### 2. Block Expansion

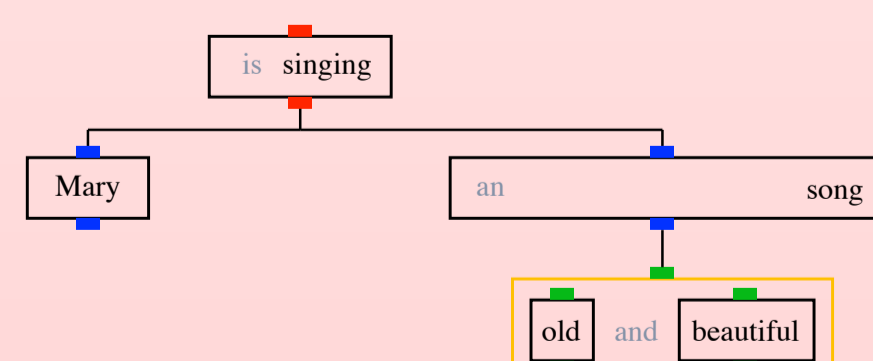
$$P_{BEM}(S) = \prod_{B \in \text{blocks}(S)} P(\text{elements}(B)|\text{derivedCat}(B))$$



Elements(B)	Derived Cat
VER	VER
NOU	NOU
an NOU	NOU
ADJ and ADJ	ADJ
ADJ	ADJ
ADJ	ADJ

#### 3. Words Filling

$$P_{WFM}(S) = \prod_{B \in \text{standardBlocks}(S)} P(\text{cw}(B)|\text{cw}(\text{parent}(B)), \text{cats}(B), \text{fw}(B), \text{dir}(B))$$



cw(B)	cw(parent(B))	cats(B)	fw(B)	direction
VBG	VER	VER	is	-
NNP	NOU	NOU	-	Left
NN	NOU	NOU	an	Right
JJ	ADJ	ADJ	old	Left
JJ	ADJ	ADJ	and	Left
JJ	ADJ	ADJ	beautiful	Left

#### 4. PoS-Tagging & Chunking

$$P_{T\&C}(S) = \prod_{i=1}^{i=\text{sentenceLength}} P(\text{tag}(i)|\text{word}_i, \text{word}_{i-1}, \text{word}_{i-2})$$

PoS	word <sub>i</sub>	pos <sub>-1</sub>	pos <sub>-2</sub>	Chunk	pos	pos <sub>-1</sub>	pos <sub>-2</sub>
NNP	Mary	-	-	N	NNP	-	-
AUX	is	NNP	-	N	AUX	NNP	-
VBG	singing	AUX	NNP	I	VBG	AUX	NNP
DT	an	VBG	AUX	N	DT	VBG	AUX
JJ	old	DT	VBG	-N	JJ	DT	VBG
JJ	and	-	-	-	-	-	-
JJ	beautiful	-	-	-	-	-	-
NN	song	-	-	-	-	-	-

$$P(S) = P_{BGM}(S) \cdot P_{BEM}(S) \cdot P_{WFM}(S) \cdot P_{T\&C}(S)$$

### Parsing through re-ranking

#### Training Phase:

- Decompose each TDS structure in the training corpus into events.
- Keep track of the frequency of each event and conditioning context (history).

#### Testing Phase:

- Get n-best PS candidates for each test sentence from another parser.
- Convert each candidate into TDS and decompose it into events.
- Compute the probability of each candidate, and output the one with maximum probability.

### Results

Treebank: Penn WSJ

Training/Test: sec 02-21 / sec 22

n-best candidates: Charniak's Max-Ent parser

#### New Evaluation metrics

- **Block Detection Score (BDS)**: the accuracy of detecting the correct boundaries of the blocks in the structure.
- **Block Attachment Score (BAS)**: the accuracy of detecting the correct governing block of each block in the structure.
- **Junction Detection Score (JDS)**: the accuracy of detecting the correct list of content-words composing each junction block in the structure.

	F-Score	UAS	BDS	BAS	JDS
Charniak (n = 1)	89.4	92.5	95.0	89.5	77.6
PCFG-reranker (n = 5)	89.0	92.4	<b>95.1</b>	89.2	77.5
PCFG-reranker (n = 1000)	83.5	88.4	92.9	83.6	71.8
TDS-reranker (n = 5)	<b>89.6</b>	92.5	<b>95.2</b>	89.5	77.6
TDS-reranker (n = 10)	89.1	92.3	95.1	89.0	76.7
TDS-reranker (n = 100)	86.9	91.0	94.4	87.0	73.1
TDS-reranker (n = 1000)	84.8	89.3	93.5	84.9	69.7

