

Towards the Cross-Roads of MWE Identification and Tree Correction [WG 4-2-3]

Agata Savary

Université François-Rabelais Tours, Blois, France



Challenge: variability of MWEs

- ▶ orthographic
 - ▷ to see the color of sb's money → to see the colour . . . ,
- ▶ morphological
 - ▷ image converters, image conversion → image converter,
- ▶ syntactic
 - ▷ the beans have been spilled → to spill the beans,
- ▶ lexical semantic
 - ▷ (FR) se **fourrer** le doigt dans l'oeil → se **mettre** le doigt dans l'oeil
(lit.) 'to put one's finger in one's eye' = 'to cherish illusions'

MWE identification in a syntax tree

- ▶ For a **contiguous** MWE:
 - ▷ generate all possible variants [3],
 - ▷ match them against the leaves of the syntax tree.
- ▶ For a **non-contiguous** MWE:
 - ▷ all instantiations correspond to a possibly infinite set of syntactic subtrees (formally: a **tree language**).

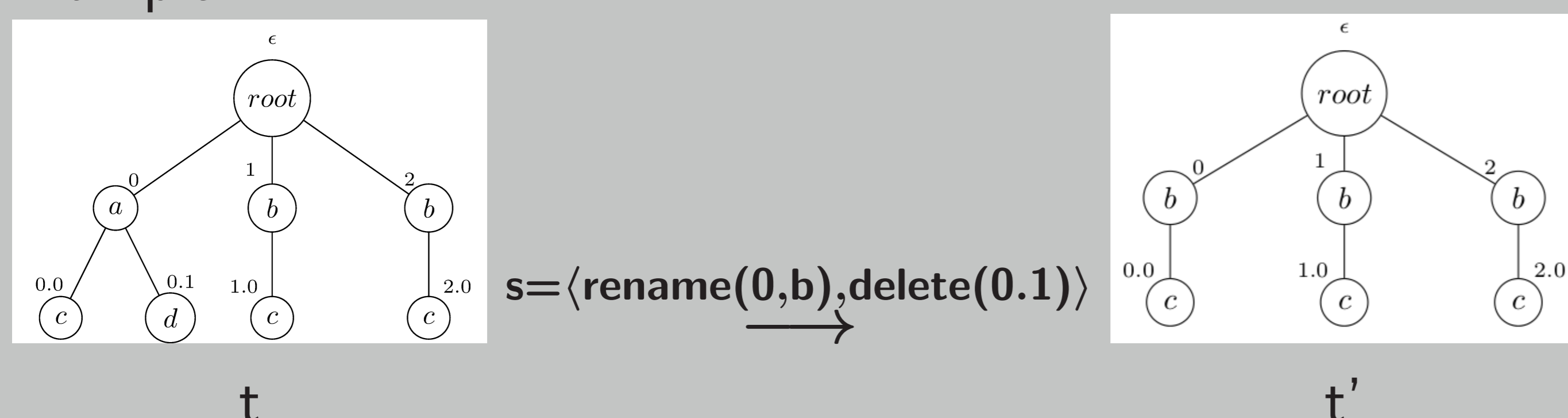
Tree-to-tree distance

- ▶ Elementary **edit operations** with **costs**, e.g.
 - ▷ relabeling a node,
 - ▷ inserting a leaf,
 - ▷ deleting a leaf (all of cost 1).
- ▶ **Edit sequences** transforming one tree into another:

$$t \xrightarrow{\text{seq}} t'$$
- ▶ **Edit distance** between trees t and t' – minimal cost of all edit sequences which transform t into t'

$$\text{dist}(t, t') = \min_{t \xrightarrow{\text{seq}} t'} \text{cost}(\text{seq})$$

▶ Example:



$$\text{dist}(t, t') = \text{cost}(s) = 2$$

Tree-to-language distance

- ▶ **Distance** between a tree t and a **tree language** L – minimal distance between t and any tree in L :

$$\text{DIST}(t, L) = \min_{t' \in L} \{\text{dist}(t, t')\}$$

Tree-to-language correction [2]

- ▶ Input:
 - ▷ tree t ,
 - ▷ tree language L ,
 - ▷ non-negative threshold th .
- ▶ Output – all trees in L whose distance from t is no higher than th

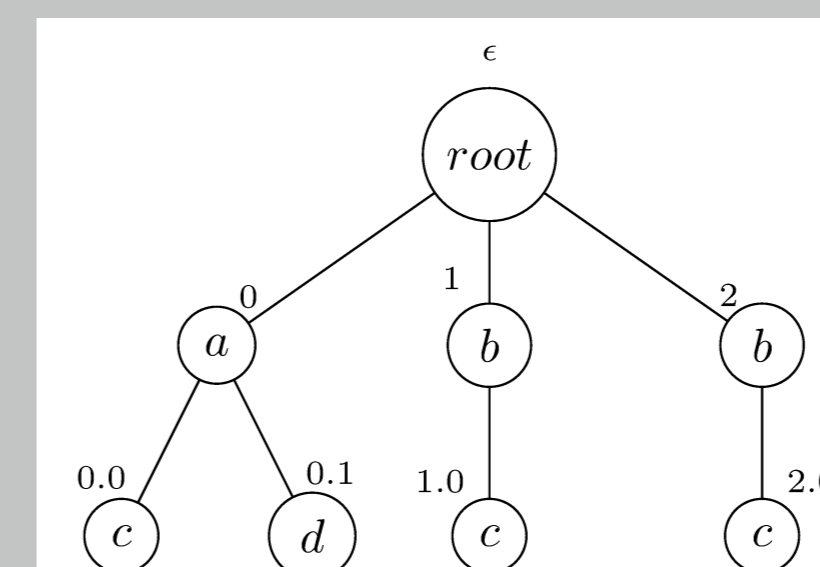
$$\text{corr}(t, L, th) = \{t' \in L : \text{dist}(t, t') \leq th\}$$

Tree-to-language correction – example

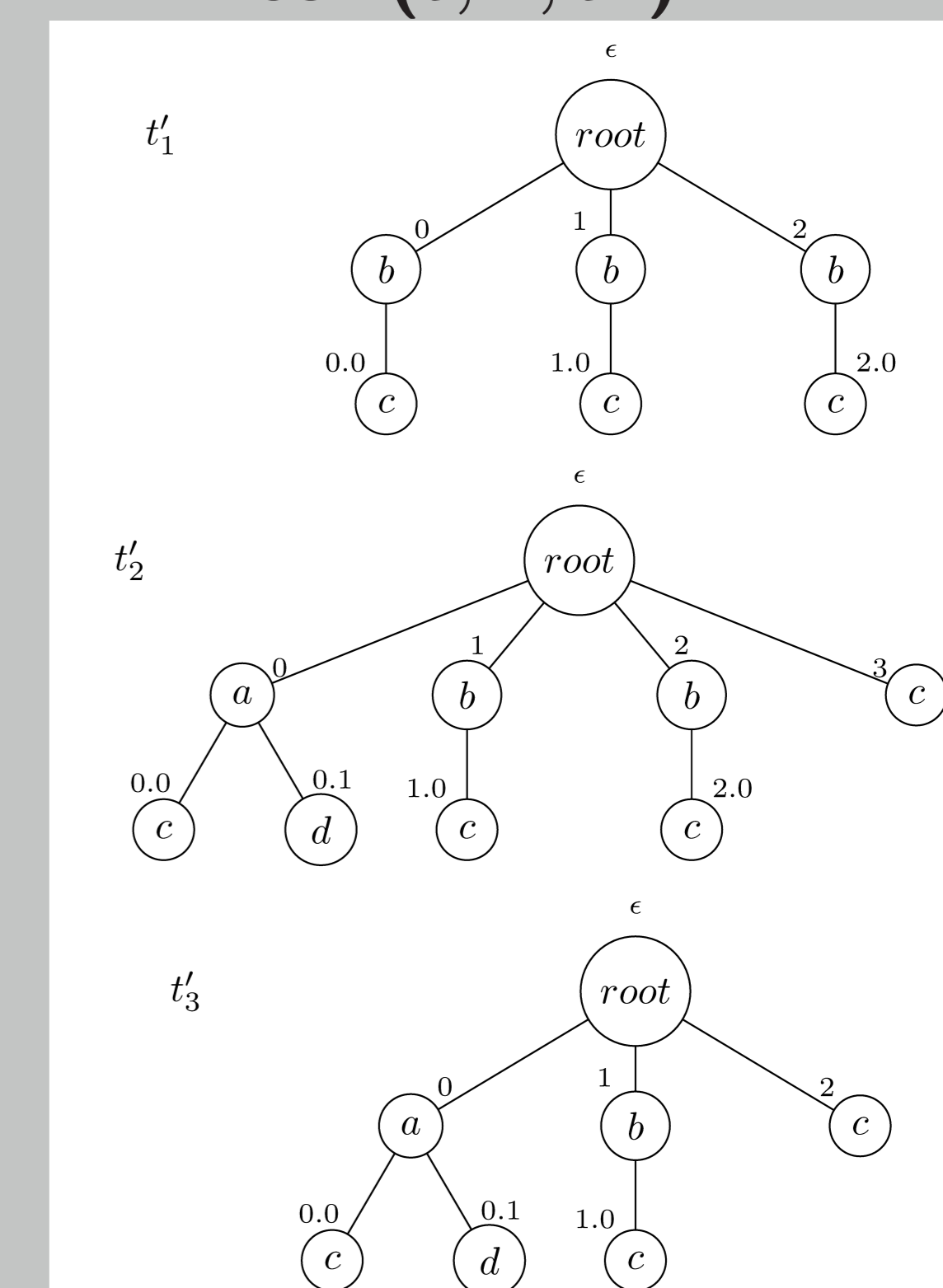
$th = 2$

$\text{DTD}(L) = \{\text{root} \rightarrow b^* | ab^*c,$
 $a \rightarrow cd,$
 $b \rightarrow c,$
 $c \rightarrow \epsilon,$
 $d \rightarrow \epsilon\}$

$t \notin L =$



$\text{corr}(t, L, th) =$



MWE identification as a tree-to-tree correction problem

- ▶ elementary operations on trees – e.g. in LTAG [1]:
 - ▷ substitution and adjunction – cost 0,
 - ▷ inserting or deleting a subtree t' at a syntactically non-allowed position – cost equal to the size of t' .
- ▶ MWE = a tree t (or a family of trees)
- ▶ occurrence of a MWE in a syntax tree = subtree t' ,
- ▶ MWE identification = finding the distance between t and t' .

MWE identification as a tree-to-language correction problem

- ▶ elementary operations on trees:
 - ▷ syntactically incorrect transformations – non-negative costs.
- ▶ MWE = tree language L (possibly infinite set of trees) – e.g.
 - ▷ *at once*: $[[[at]_{\text{Prep}}[once]_{\text{Adv}}]_{\text{AdvP}}[s]]_s$
 - ▷ $L_{\text{at once}}$ – set of all trees that result from its auxiliary tree by its adjunction to any other tree.
- ▶ occurrence of a MWE in a syntax tree = subtree t ,
- ▶ MWE identification = correcting t with respect to L under a given threshold th .

Applications

- ▶ post-annotating MWEs in treebanks,
- ▶ detecting MWEs in a post-parsing stage,
- ▶ when $th > 0$:
 - ▷ processing noisy data (spontaneous speech, social networks),
 - ▷ detecting errors in corpus annotation, grammar, or MWE lexicon.

Bibliography

- [1] Anne Abeillé and Yves Schabes.
Parsing idioms in lexicalized tags.
 In *EACL'89, Manchester*, pages 1–9, 1989.
- [2] Joshua Amavi, Béatrice Bouchou, and Agata Savary.
On Correcting XML Documents with Respect to a Schema.
The Computer Journal, 2013.
- [3] Agata Savary.
Multiflex: A Multilingual Finite-State Tool for Multi-Word Units.
 volume 5642 of *Lecture Notes in Computer Science*, pages 237–240.
 Springer, 2009.