

Model-Checking Parse Trees

A. Boral S. Schmitz

CMI, Chennai LSV, ENS Cachan

LICS 2013, June 26th 2013

MODELING SYNTAX (WITH TREES)

(E.G. PULLUM AND SCHOLZ, 2001)

generative syntax context-free grammars, ...

model-theoretic syntax MSO, PDL_{tree}, ...

mixed here for grammar engineering

Outline

Model Checking Parse Trees

Computational Linguistics

Programming Language Design

Complexity Results

MODELING SYNTAX (WITH TREES)

(E.G. PULLUM AND SCHOLZ, 2001)

generative syntax context-free grammars, ...

model-theoretic syntax MSO, PDL_{tree}, ...

mixed here for grammar engineering

Outline

Model Checking Parse Trees

Computational Linguistics

Programming Language Design

Complexity Results

RECOGNITION ACROSS FORMALISMS

CFG local specification PTIME-c.

PDL_{tree} long-distance navigation EXPTime-c.

MSO global specification TOWER-c.

PARSE FOREST MODEL-CHECKING

input a context free grammar G , a word w ,
and a PDL_{tree} formula φ ,

question does there exists a tree $t \in L_{G,w}$ s.t.
 $t \models \varphi$?

PDL ON TREES

(AFANASIEV et al., 2005)

- ▶ finite, ordered, AP-labeled trees
- ▶ two relations: **child** ' \downarrow ', **right-sibling** ' \rightarrow '
- ▶ syntax of PDL_{tree} (aka Regular XPath):

$$\varphi ::= p \mid \top \mid \neg\varphi \mid \varphi \wedge \varphi \mid \langle \pi \rangle \varphi \quad (\text{node formulæ})$$

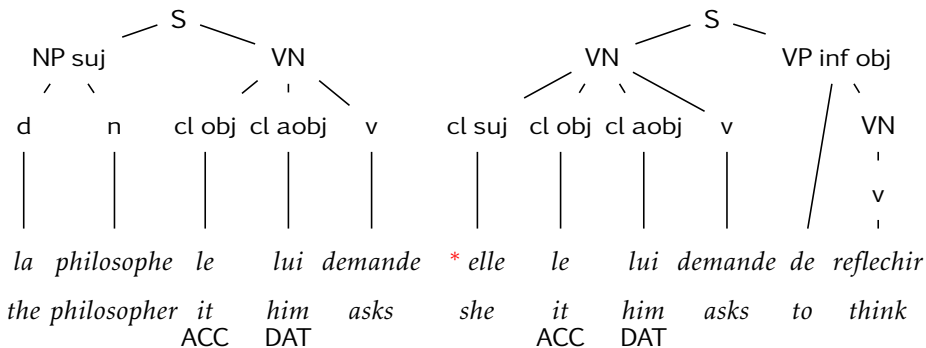
$$\pi ::= \downarrow \mid \rightarrow \mid \pi; \pi \mid \pi + \pi \mid \pi^* \mid \pi^{-1} \mid \varphi? \\ (\text{path formulæ})$$

LONG-DISTANCE DEPENDENCIES

- ▶ **local** syntax using a CFG
- ▶ **long-distance** constraints enforced through a PDL_{tree} formula

EXAMPLE: FRENCH CLITICS

$$S \rightarrow \text{NP}_{\text{suj}}? \text{VN} \text{VP}_{\text{infobj}}? \text{PP}_{\text{aobj}}?$$

$$\text{VN} \rightarrow \text{cl}_{\text{suj}}? \text{cl}_{\text{obj}}? \text{cl}_{\text{aobj}}? v$$


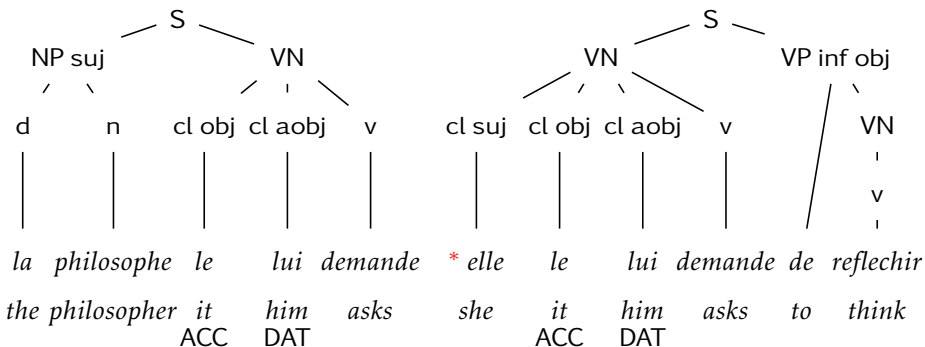
EXAMPLE: FRENCH CLITICS

$[\downarrow^*]demande \Rightarrow \langle (\uparrow; \uparrow; \rightarrow^+) + (\uparrow; \leftarrow^+; cl?) \rangle obj$ (at least one object)

$\wedge \langle (\uparrow; \uparrow; \leftarrow) + (\uparrow; \leftarrow^+; cl?) \rangle suj$ (at least one subject)

$\wedge \bigwedge_{f \in \{suj, obj, aobj\}} \langle \uparrow; \leftarrow^+; cl? \rangle f \Rightarrow \neg \langle \uparrow; \uparrow; (\leftarrow + \rightarrow^+) \rangle f$

(clitic arguments forbid the matching canonical arguments)



DISCUSSION

grammar ▶ **ϵ -free**: $A \rightarrow \epsilon$ forbidden

- ▶ **acyclic**: $A \Rightarrow^+ A$ forbidden

formula ▶ here PDL_{core} (aka Core XPath)

- ▶ in the literature: PDL_{cp} (aka Conditional XPath, see Palm, 1999; Lai and Bird, 2010)

model ▶ richer structures than trees?

- ▶ **trace** can be emulated by a formula

AMBIGUITY FILTERING

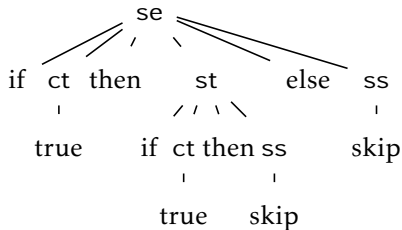
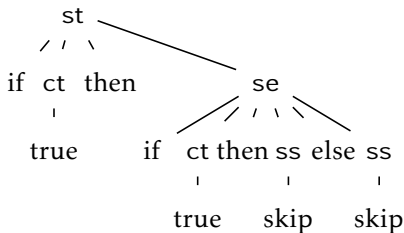
(E.G. KLINT AND VISSER, 1994)

- ▶ syntax defined through a CFG
- ▶ exclude spurious parse trees, here thanks to a formula

EXAMPLE: DANGLING ELSE

$$S \rightarrow st(\text{if } C \text{ then } S) \mid se(\text{if } C \text{ then } S \text{ else } S)$$

$$\mid sw(\text{while } C \text{ } S) \mid ss(\text{skip})$$

$$C \rightarrow ct(\text{true}) \mid cf(\text{false})$$


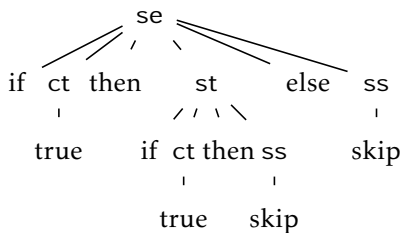
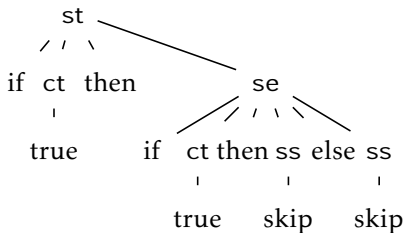
EXAMPLE: DANGLING ELSE

$$\text{first} \stackrel{\text{def}}{=} \neg \langle \leftarrow \rangle \top$$

$$\text{last} \stackrel{\text{def}}{=} \neg \langle \rightarrow \rangle \top$$

$$\prec \stackrel{\text{def}}{=} (\text{last?}; \uparrow)^*; \rightarrow; (\downarrow; \text{first?})^*$$

$$\varphi = \neg \langle \downarrow^* \rangle (\text{st} \wedge \langle \prec \rangle \text{else})$$



DISCUSSION

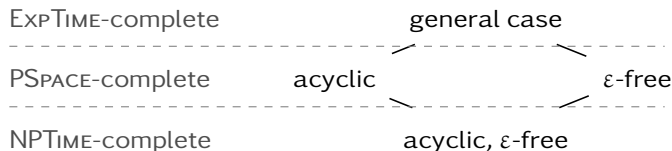
grammar ▶ acyclic grammars

formula ▶ here in PDL_{cp} (but also in XPath 1.0)

model ▶ “static” filtering: can be compiled into the grammar

- ▶ handling layout-sensitive syntax (aka offside rules)? (c.f. Adams, 2013)

COMPLEXITY SUMMARY



- ▶ borrow from similar problems in XML theory by Benedikt et al. (2008)
- ▶ lower bounds already hold for the $\text{PDL}_{\text{core}}[\downarrow]$ fragment, often for fixed G and/or w .

PARSE FOREST $L_{G,w}$

Theorem (Bar-Hillel et al., 1961)

Let G be a CFG and w a string. Then we can construct in polynomial time a tree automaton $\mathcal{A}_{G,w}$ of polynomial size, s.t. $L(\mathcal{A}_{G,w}) = L_{G,w}$.

- ▶ intuition: the states of $\mathcal{A}_{G,w}$ are triples (i, A, j) where $0 \leq i \leq j \leq |w|$ and $w_i \cdots w_j \in L_G(A)$.
- ▶ cardinality of $L_{G,w}$: potentially



PARSE FOREST $L_{G,w}$

Theorem (Bar-Hillel et al., 1961)

Let G be a CFG and w a string. Then we can construct in polynomial time a tree automaton $\mathcal{A}_{G,w}$ of polynomial size, s.t. $L(\mathcal{A}_{G,w}) = L_{G,w}$.

- ▶ intuition: the states of $\mathcal{A}_{G,w}$ are triples (i, A, j) where $0 \leq i \leq j \leq |w|$ and $w_i \cdots w_j \in L_G(A)$.
- ▶ cardinality of $L_{G,w}$: potentially

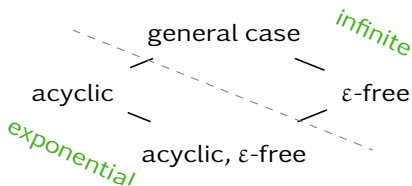


PARSE FOREST $L_{G,w}$

Theorem (Bar-Hillel et al., 1961)

Let G be a CFG and w a string. Then we can construct in polynomial time a tree automaton $\mathcal{A}_{G,w}$ of polynomial size, s.t. $L(\mathcal{A}_{G,w}) = L_{G,w}$.

- ▶ intuition: the states of $\mathcal{A}_{G,w}$ are triples (i, A, j) where $0 \leq i \leq j \leq |w|$ and $w_i \cdots w_j \in L_G(A)$.
- ▶ cardinality of $L_{G,w}$: potentially



ZOOM ON PSPACE RESULTS

- ▶ build a nonrecursive DTD D from G and w
 - ▶ **DTD**: “extended” CFG with regular right parts
 - ▶ thus **unranked** trees
 - ▶ **nonrecursive**: $A \Rightarrow^+ \alpha A \beta$ forbidden
- ▶ build a formula φ'

Theorem (Benedikt et al., 2008)

Satisfiability of a PDL_{core} formula φ' in presence of a nonrecursive DTD D is in PSPACE.

(Reduces to emptiness of loop-free alternating two-ways word automata on an XML encoding (Serre, 2006).)

ZOOM ON PSPACE RESULTS

- ▶ build a nonrecursive DTD D from G and w
 - ▶ **DTD**: “extended” CFG with regular right parts
 - ▶ thus **unranked** trees
 - ▶ **nonrecursive**: $A \Rightarrow^+ \alpha A \beta$ forbidden
- ▶ build a formula φ'

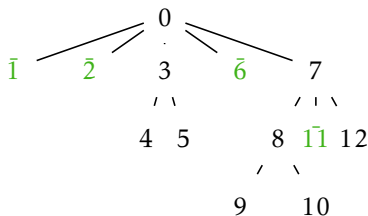
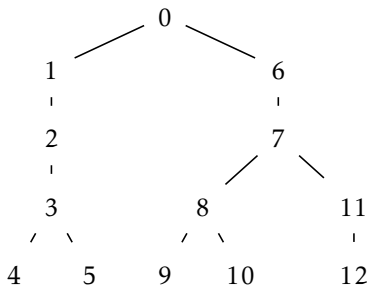
Theorem

*Satisfiability of a **PDL_{tree}** formula φ' in presence of a nonrecursive DTD D is in PSPACE.*

(Reduces to emptiness of alternating parity two-ways word automata on an XML encoding (Serre, 2006).)

REDUCTION IN THE ϵ -FREE CASE

- ▶ $L_{G,w}$ can be infinite
- ▶ **chains** of transitions in $\mathcal{A}_{G,w}$: $q_1 \Rightarrow q_2 \Rightarrow \dots \Rightarrow q_n$
- ▶ **maximal** chain if q_n binary or a leaf
- ▶ “rotate” chains



REDUCTION IN THE ε -FREE CASE (CONT'D)

- ▶ $\text{chains}(q)$ the set of maximal chains starting from q forms a regular substitution
- ▶ D obtained from $\mathcal{A}_{G,w}$ by applying $\text{chains}(q)$
- ▶ φ has to be interpreted on the “rotated” trees

Proposition

PFMC with ε -free grammars is in PSPACE.

Corollary

PDL_{tree} recognition with ε -free models is in PSPACE.

CONCLUDING REMARKS

- ▶ **mixed** approach to syntax:
 - ▶ bulk work with a CFG
 - ▶ fine tuning with a PDL_{tree} (or MSO, or...) formula
- ▶ new applications for model-checking techniques
 - ▶ different practical restrictions
 - ▶ leading to improved complexities
- ▶ exploit data logics?

REFERENCES

- Adams, M.D., 2013. Principled parsing for indentation-sensitive languages: revisiting Landin's offside rule. In *POPL 2013*, pages 511–522. doi:10.1145/2429069.2429129.
- Afanasiev, L., Blackburn, P., Dimitriou, I., Gaiffe, B., Goris, E., Marx, M., and de Rijke, M., 2005. PDL for ordered trees. *J. Appl. Non-Classical Log.*, 15(2):115–135. doi:10.3166/jancl.15.115-135.
- Bar-Hillel, Y., Perles, M., and Shamir, E., 1961. On formal properties of simple phrase-structure grammars. *Z. Phonetik Sprachwiss. Kommunik.*, 14:143–172.
- Benedikt, M., Fan, W., and Geerts, F., 2008. XPath satisfiability in the presence of DTDs. *Journal of the ACM*, 55(2:8). doi:10.1145/1346330.1346333.
- Klint, P. and Visser, E., 1994. Using filters for the disambiguation of context-free grammars. In Pighizzini, G. and San Pietro, P., editors, *ASMIACS Workshop on Parsing Theory*, Technical Report 126-1994, pages 89–100. Università di Milano.
- Lai, C. and Bird, S., 2010. Querying linguistic trees. *J. Logic Lang. Inform.*, 19(1):53–73. doi:10.1007/s10849-009-9086-9.
- Palm, A., 1999. Propositional tense logic of finite trees. In *MOL 6*.
- Pullum, G.K. and Scholz, B.C., 2001. On the distinction between model-theoretic and generative-enumerative syntactic frameworks. In *LACL 2001*, volume 2099 of *LNCS*, pages 17–43. Springer. doi:10.1007/3-540-48199-0_2.
- Serre, O., 2006. Parity games played on transition graphs of one-counter processes. In Aceto, L. and Ingólfssdóttir, A., editors, *FoSSaCS 2006*, volume 3921 of *LNCS*, pages 337–351. Springer. doi:10.1007/11690634_23.