

PDL for Ordered Trees

Loredana Afanasiev, Patrick Blackburn, Ioanna Dimitriou,
Bertrand Gaiffe, Evan Goris, Maarten Marx, Maarten de Rijke

May 17, 2004

Abstract

This paper is about a special version of PDL, designed by Marcus Kracht for reasoning about sibling ordered trees. It has four basic programs corresponding to the child, parent, left- and right-sibling relations in such trees. The original motivation for this language is rooted in the field of model-theoretic syntax. Motivated by recent developments in the area of semistructured data, and, especially, in the field of query languages for XML (eXtensible Markup Language) documents, we revisit the language. This renewed interest comes with a special focus on complexity and expressivity aspects of the language, aspects that have so far largely been ignored. We survey and derive complexity results, and spend most of the paper on the most important open question concerning the language: what is its expressive power? We approach this question from two angles: Which first order properties can be expressed? And which second order properties? While we are still some way from definitive answers to these questions, we discuss two first order fragments of the PDL language for ordered trees, and show how the language can be used to express some typical (second order) problems, like the boolean circuit and the frontier problem.

1 Introduction

The purpose of this paper is to revive interest in a version of PDL proposed by Marcus Kracht [18, 19]. This version, here called PDL_{tree} , is specially designed for models which are sibling ordered trees. Such models are of interest in at least two research communities: linguistics, in particular the field of model-theoretic syntax, and computer science, in particular those working with the world wide web, semi structured data and XML (eXtensible Markup Language).

Model-theoretic syntax is an uncompromisingly declarative approach to natural language syntax: grammatical theories are logical theories, and grammatical structures are their models. These models consists of parse trees, i.e., node labelled, sibling ordered finite trees. Perhaps the best known work in this tradition is that of James Rogers (for example [27]) in which

grammatical theories are stated in monadic second-order logic. However other authors (in particular [4, 18, 19, 25]) use various kinds of *modal logic* (in essence, variable free formalisms for describing relational structures) to specify grammatical constraints. Palm [25] contains some interesting linguistic examples and is a good introduction to (and motivation for) this approach.

The World Wide Web is a freely evolving, ever-changing collection of data with flexible structure. The Web's nature escapes the conventional database scenario of manipulating data: data on the Web simply do not comply with the strict schemas used for conventional databases. Web data such as home pages, news sites, pages on commercial sites, usually enjoy some amount of structure, but that is not strictly enforced, and there are no uniformly adopted standards, not even for simple bits of information such as addresses. Hence, data on the Web is essentially semi-structured [1]. In search for suitable models for semi-structured data, the World Wide Web Consortium proposed the eXtensible Markup Language (XML) [6]. XML is a standard for textual representation of semi-structured data and was designed to describe any type of textual information. It looks like a flexible variant of HTML, allowing for the mark-up of data with information about its content rather than its presentation. The logical abstraction of an XML document (the so called DOM) is a finite node labelled ordered tree.

Motivated by the renewed need for clean, well-understood declarative tree description formalisms brought about by the developments in semistructured data outlined above, we want to revive interest in the special variant of PDL developed for sibling ordered trees. We focus on complexity and expressivity aspects of the language. Section 2 introduces the language. Section 3 discusses complexity, and in Section 4 and Section 5 we address expressivity issues. Section 4 is devoted to the expressiveness of the language in terms of first order properties; we discuss the first order fragment of PDL_{tree} , recall some known results, and show the language in action by expressing the until modality over the document order relation.

Several proposals for languages that are complete for unary monadic second order logic have been made, but none of these is as simple and easy to learn as PDL_{tree} . So the most pressing issue seems to be to determine the exact expressive power of PDL_{tree} compared to unary monadic second order logic. This remains an open problem, but to improve our understanding of PDL_{tree} 's expressive power, we adopt a well-known strategy by examining a number of 'typical' second order problems and properties. Specifically, in Section 5 we show how we can express the boolean circuit and the frontier problem, and we discuss infinity axioms. These examples suggest that PDL_{tree} is expressive enough to encode natural hard second order problems. Boolean circuits is one of the main problems used to show that a logic is weaker than MSO. The frontier problem is a typical linguistical application. Being able to express finiteness shows a certain robustness of the language.

We conclude in Section 6.

2 PDL for ordered trees

We recall the definition of PDL_{tree} from [18, 19]. PDL_{tree} is a propositional modal language identical to Propositional Dynamic Logic (PDL) [16] over four basic programs `left`, `right`, `up` and `down`, which explore the left-sister, right-sister, mother-of and daughter-of relations. Recall that PDL has two sorts of expressions: programs and propositions. We suppose we have fixed a non-empty, finite or countably infinite, set of atomic symbols A whose elements are typically denoted by p . PDL_{tree} 's syntax is as follows, writing π for programs and ϕ for propositions:

$$\begin{aligned}\pi &::= \text{left} \mid \text{right} \mid \text{up} \mid \text{down} \mid \pi; \pi \mid \pi \cup \pi \mid \pi^* \mid ?\phi \\ \phi &::= p \mid \top \mid \neg\phi \mid \phi \wedge \phi \mid \langle \pi \rangle \phi.\end{aligned}$$

We sometimes write $\text{PDL}_{\text{tree}}(A)$ to emphasize the dependence on A . We employ the usual boolean abbreviations and use $[\pi]\phi$ for $\neg\langle \pi \rangle\neg\phi$.

We interpret $\text{PDL}_{\text{tree}}(A)$ on *finite ordered trees* whose nodes are *labeled* with symbols drawn from A . We assume that the reader is familiar with finite trees and such concepts as ‘daughter-of’, ‘mother-of’, ‘sister-of’, ‘root-node’, ‘terminal-node’, and so on. If a node has no sister to its immediate right we call it a last node, and if it has no sister to its immediate left we call it a first node. Note that the root node is both first and last. The root node will always be called *root*. A labelling of a finite tree associates a subset of A with each tree node.

Formally, we present finite ordered trees¹ (tree for short) as tuples $\mathbf{T} = (T, R_{\text{right}}, R_{\text{down}})$. Here T is the set of tree nodes and R_{right} and R_{down} are the right-sister and daughter-of relations respectively. A pair $\mathfrak{M} = (\mathbf{T}, V)$, where \mathbf{T} is a finite tree and $V : A \rightarrow \text{Pow}(T)$, is called a *model*, and we say that V is a *labelling function* or a *valuation*. Given a model \mathfrak{M} , we simultaneously define a set of relations on $T \times T$ and the interpretation of the language $\text{PDL}_{\text{tree}}(A)$ on \mathfrak{M} :

$$\begin{aligned}R_{\text{up}} &= R_{\text{down}}^{-1} & R_{\pi \cup \pi'} &= R_{\pi} \cup R_{\pi'} \\ R_{\text{left}} &= R_{\text{right}}^{-1} & R_{\pi; \pi'} &= R_{\pi} \circ R_{\pi'} \\ R_{\pi^*} &= R_{\pi}^* & R_{?\phi} &= \{(t, t) \mid \mathfrak{M}, t \models \phi\}.\end{aligned}$$

¹A sibling ordered tree is a structure isomorphic to $(N, R_{\downarrow}, R_{\rightarrow})$ where N is a set of finite sequences of natural numbers closed under taking initial segments, and for any sequence s , if $s \cdot k \in N$, then either $k = 0$ or $s \cdot k - 1 \in N$. For $n, n' \in N$, $nR_{\downarrow}n'$ holds iff $n' = n \cdot k$ for k a natural number; $nR_{\rightarrow}n'$ holds iff $n = s \cdot k$ and $n' = s \cdot k + 1$.

$$\begin{aligned}
\mathfrak{M}, t \models p & \text{ iff } t \in V(p), \text{ for all } p \in A \\
\mathfrak{M}, t \models \top & \text{ iff } t \in T \\
\mathfrak{M}, t \models \neg\phi & \text{ iff } \mathfrak{M}, t \not\models \phi \\
\mathfrak{M}, t \models \phi \wedge \psi & \text{ iff } \mathfrak{M}, t \models \phi \text{ and } \mathfrak{M}, t \models \psi \\
\mathfrak{M}, t \models \langle \pi \rangle \phi & \text{ iff } \exists t' (tR_\pi t' \text{ and } \mathfrak{M}, t' \models \phi).
\end{aligned}$$

If $\mathfrak{M}, t \models \phi$, then we say ϕ is *satisfied* in \mathfrak{M} at t . For any formula ϕ , if there is a model \mathfrak{M} such that $\mathfrak{M}, \text{root} \models \phi$, then we say that ϕ is *satisfiable*. For Γ a set of formulas, and ϕ a formula, we say that ϕ is a consequence of Γ (denoted by $\Gamma \models \phi$) if for every model in which Γ is satisfied at every node, ϕ is also satisfied at every node.

Note that we could have defined PDL_{tree} by taking **down** and **right** as the sole primitive programs and closing the set of programs under converses. As converse commutes with all program operators, these two definitions are equally expressive.

Let's consider some examples: if universally true, (1) says that every a node has a b and a c daughter, in that order, and no other daughters; and (2) says that every a node has a b first daughter followed by some number of c daughters, and no other daughters.

$$a \rightarrow \langle \text{down} \rangle (\neg \langle \text{left} \rangle \top \wedge b \wedge \langle \text{right} \rangle (c \wedge \neg \langle \text{right} \rangle \top)) \quad (1)$$

$$a \rightarrow \langle \text{down} \rangle (\neg \langle \text{left} \rangle \top \wedge b \wedge \langle \langle \text{right}; ?c \rangle^* \rangle \neg \langle \text{right} \rangle \top). \quad (2)$$

Now consider (3). This projects a label p down to some leaf node:

$$\langle \langle ?p; \text{down} \rangle^* \rangle (p \wedge \neg \langle \text{down} \rangle \top) \quad (3)$$

That is, whenever this formula is satisfied in some model at some point t , there will be a path from t to some leaf node l such that every node on the path is marked p . We end the short examples with a list of useful abbreviations:

abbreviation	of
root	$\neg \langle \text{up} \rangle \top$
leaf	$\neg \langle \text{down} \rangle \top$
first	$\neg \langle \text{left} \rangle \top$
last	$\neg \langle \text{right} \rangle \top$
π^+	$\pi; \pi^*$.

3 Complexity

We now look at the complexity of the PDL_{tree} consequence problem: how difficult is it to decide whether, on finite ordered trees, $\Gamma \models \chi$, for finite Γ .

Decidability of this problem follows from the interpretation of PDL_{tree} into $L_{K,P}^2$ [27]. (The decidability of the satisfiability problem for $L_{K,P}^2$ follows, in turn, via an interpretation into Rabin's $S\omega S$.) But although this reduction yields decidability, it only gives us a non-elementary decision procedure. So what is the complexity of the consequence problem?

Recall that EXPTIME is the class of all problems solvable in exponential time. A problem is solvable in exponential time if there is a deterministic exponentially time bounded Turing machine that solves it. A deterministic Turing machine is exponentially time bounded if there is a polynomial $p(n)$ such that the machine always halts after at most $2^{p(n)}$ steps, where n is the length of the input. To prove EXPTIME -completeness we have to do two things: prove an EXPTIME lower bound (that is, show that some problem instances require exponential time) and an EXPTIME upper bound (that is, give an algorithm that handles any problem instance in exponential time). Let's first deal with the lower bound.

Theorem 3.1 ([12, 28]). *The consequence problem for the PDL_{tree} fragment with only down is EXPTIME -hard.*

Proof. This is an immediate corollary of the analysis of the lower bound result for PDL given by [28], based on the work of [12]. She notes that the following fragment of PDL is EXPTIME -hard: formulas of the form $\psi \wedge [a^*]\theta$ (where ψ and θ contain only the atomic program a and no embedded modalities) that are satisfiable at the root of a finite binary tree. Identifying the program a with \downarrow , the result follows (because $[\text{down}^*]\theta \wedge \psi$ is satisfiable at the root of a finite tree iff $\theta \not\models \text{root} \rightarrow \neg\psi$). \dashv

For full PDL this bound is optimal. There is even a stronger result: every satisfiable PDL formula ϕ can be satisfied on a model with size exponential in the length of ϕ . However with tree-based models there is no hope for such a result for it is easy to show that:

For every natural number n , there exists a satisfiable formula of size $\mathcal{O}(n^2)$ in the language with only \downarrow and \downarrow^ which can only be satisfied on at least binary branching trees of depth at least 2^n .*

A formula containing most of the requirements to force such a deep branch is given in Proposition 6.51 of [2]. To this formula we only have to add the conjunct $[\text{down}^*](\langle\downarrow\rangle p \wedge \langle\downarrow\rangle \neg p)$ for some new variable p to enforce binary branching. Note that the size of such a model is double exponential in the size of the formula. This means that a decision algorithm which tries to construct a tree model must use at least exponential space, as it will need to keep a whole branch in memory.

So we're going to have to think more carefully about the upper bound. One way to proceed is to take a cue from the completeness proof for a related language in [4]. Instead of constructing a model it is possible to

design an algorithm which searches for a “good” set of labellings of the nodes of a model. Label sets consist of subformulas of the formula ϕ whose satisfiability is to be decided. From a good set of labels we can construct a labeled tree model which satisfies ϕ . The number of labels is bound by an exponential in the number of subformulas of ϕ , and the search for a good set of labels among the possible ones can be implemented in time polynomial in the number of possible labels using the technique of elimination of Hintikka sets developed in [26]. A direct proof using this technique was given in [3] for the language \mathcal{L}_{cp} (see Section 4). Unfortunately, the technique cannot be straightforwardly applied to PDL_{tree} . Here we show how an old result of Vardi and Wolper [29] on deterministic PDL with converse yields the desired upper bound.

Theorem 3.2. *The PDL_{tree} consequence problem is in EXPTIME.*

Proof. First note that $\gamma_1, \dots, \gamma_n \models \chi$ if and only if $\models \text{root} \rightarrow ([\text{down}^*](\gamma_1 \wedge \dots \wedge \gamma_n) \rightarrow \chi)$. Thus we need only decide satisfiability of PDL_{tree} formulas at the root of finite trees.

Now consider the language \mathcal{L}_2 , the modal language with only the two programs $\{\downarrow_1, \downarrow_2\}$ and their inverses $\{\uparrow_1, \uparrow_2\}$. \mathcal{L}_2 is interpreted on finite at most *binary-branching* trees, with \downarrow_1 and \downarrow_2 interpreted by the first and second daughter relation, respectively. We will effectively reduce PDL_{tree} satisfiability to \mathcal{L}_2 satisfiability. \mathcal{L}_2 is a fragment of deterministic PDL with converse. [29] shows that the satisfiability problem for this latter language is decidable in EXPTIME over the class of all models. This is done by constructing for each formula ϕ a tree automaton A_ϕ which accepts exactly all tree models in which ϕ is satisfied. Thus deciding satisfiability of ϕ reduces to checking emptiness of A_ϕ . The last check can be done in time polynomial in the size of A_ϕ . As the size of A_ϕ is exponential in the length of ϕ , this yields the exponential time decision procedure.

But we want satisfiability on *finite* trees. This is easy to cope with in an automata-theoretic framework: construct an automaton $A_{\text{fin_tree}}$, which accepts only finite binary trees, and check emptiness of $A_\phi \cap A_{\text{fin_tree}}$. The size of $A_{\text{fin_tree}}$ does not depend on ϕ , so this problem is still in EXPTIME.

The reduction from PDL_{tree} to \mathcal{L}_2 formulas is very simple: replace the PDL_{tree} programs `down`, `up`, `right`, `left` by the \mathcal{L}_2 programs

$$\downarrow_1; \downarrow_2^*, \uparrow_2^*; \uparrow_1, \downarrow_2, \uparrow_2,$$

respectively. It is straightforward to prove that this reduction preserves satisfiability, following the reduction from $S\omega S$ to $S2S$ as explained in [30]: an PDL_{tree} model $(T, R_{\text{right}}, R_{\text{down}}, V)$ is turned into an \mathcal{L}_2 model (T, R_1, R_2, V) by defining

$$R_1 = \{(x, y) \mid xR_{\text{down}}y \text{ and } y \text{ is the first daughter of } x\}$$

and $R_2 = R_{\text{right}}$. Turn an \mathcal{L}_2 model (T, R_1, R_2, V) into an PDL_{tree} model $(T, R_{\text{right}}, R_{\text{down}}, V)$ by defining $R_{\text{right}} = R_2$ and $R_{\text{down}} = R_1 \circ R_2^*$. \dashv

4 Expressivity 1: First Order Logic

Let \mathcal{L}_{FO} denote the first-order language over the signature with binary predicates $\{R_{\text{down}^+}, R_{\text{right}^+}\}$ and countably many unary predicates. \mathcal{L}_{FO} is interpreted on ordered trees in the obvious way: R_{down^+} is interpreted by the transitive closure of the daughter relation, and R_{left^+} is interpreted by the transitive closure of the left-sister relation. Note that the language is first order, even though we interpret the two primitive relations as second order relations over a more primitive relation. This is not strange, but just another perspective: we take descendant as primitive instead of the immediate daughter relation. Of course the latter is first order definable from the descendant relation.

Two other modal languages proposed in the model-theoretic syntax literature can be considered as first order fragments of PDL_{tree} . That is, they can be considered as versions of PDL_{tree} with a more limited repertoire of programs. Palm [25] even argues that for linguistic applications one *must* restrict the language to its first order fragment. As first order logic is a natural point of reference for the expressivity of languages it is useful to consider first order fragments of PDL_{tree} . We consider two, one predating and one postdating the introduction of PDL_{tree} .

The language proposed by Blackburn, Meyer-Viol and de Rijke [5], here called $\mathcal{L}_{\text{Core}}$, contains only the core machinery for describing trees: the four basic programs plus their *transitive* closures, denoted by a superscript $(\cdot)^+$. This language is precisely as expressive as the fragment of PDL_{tree} generated by the following programs:

$$\pi ::= \text{left} \mid \text{right} \mid \text{up} \mid \text{down} \mid \pi^*,$$

or equivalently by

$$\pi ::= \text{left} \mid \text{right} \mid \text{up} \mid \text{down} \mid \pi; \pi \mid \pi \cup \pi \mid ?\phi \mid a^*, \text{ for } a \text{ one of the four atomic programs.}$$

The language proposed by Palm [25], here called \mathcal{L}_{cp} (with *cp* abbreviating *conditional path*), lies between $\mathcal{L}_{\text{Core}}$ and PDL_{tree} with respect to expressive power. It can be thought of as the fragment of PDL_{tree} generated by the following programs:

$$\pi ::= \text{left} \mid \text{right} \mid \text{up} \mid \text{down} \mid \pi; ?\phi \mid \pi^*,$$

or equivalently by

$$\pi ::= \text{left} \mid \text{right} \mid \text{up} \mid \text{down} \mid \pi; \pi \mid \pi \cup \pi \mid ?\phi \mid (a; ?\phi)^*, \text{ for } a \text{ one of the four atomic programs.}$$

Note that $(a; ?\phi)^+$ and $(?\phi; a)^*$ are definable, and we will use these as abbreviations below. Both languages are easily seen to be fragments of \mathcal{L}_{FO} , the first order language for ordered trees. In fact we know exactly which fragments:

Theorem 4.1 ([20]). *The following are equivalent on ordered trees. For N a set of nodes:*

- N is definable by an \mathcal{L}_{cp} formula;
- N is definable by an \mathcal{L}_{FO} formula in one free variable.

Theorem 4.2 ([21]). *The following are equivalent on ordered trees. For N a set of nodes:*

- N is definable by an \mathcal{L}_{Core} formula;
- N is definable by an \mathcal{L}_{FO} formula in one free variable which
 1. contains at most two (free and bound) variables (possibly reused), and
 2. which may use additional atomic relations corresponding to the right-sister and daughter-of relation.

The first theorem can be seen as a generalization of Kamp’s Theorem [17] to ordered trees. The theorem was announced in [24], but the proof is hard to follow. [20] contains a proof based on Gabbay’s notion of separation [13]. The second theorem is also a generalization of a result for temporal logic on linear structures, this time due to Etessami, Vardi and Wilke [11].

We end this section by giving some insight in the expressive power of \mathcal{L}_{cp} . First note that the temporal until(ϕ, ψ) modality can be expressed, in all four directions. For the downward direction, until(ϕ, ψ) is expressed as

$$\langle\langle \mathbf{down}; ?\psi \rangle^*\rangle \langle \mathbf{down} \rangle \phi.$$

Indeed, this formula is true at a node n if and only if there exists a desendent n' of n at which ϕ is true and at all nodes strictly in between ψ is true.

In the context of XML documents, the order in which the nodes are written is an important relation, called document order. Figure 1 contains an example of an XML file, its corresponding tree model (called the DOM) and the ordering. The document order relation \ll is defined as

$$\ll \equiv \mathbf{down}^+ \cup \mathbf{up}^*; \mathbf{right}^+; \mathbf{down}^*.$$

On finite trees it makes sense to speak about the successor relation of the document order. The simple definition is $\ll \cap \overline{\ll} \circ \ll$. It can be defined also with the \mathcal{L}_{cp} programs as

$$\mathbf{down}; ?\mathbf{first} \cup ?\mathbf{leaf}; \mathbf{right} \cup (? \mathbf{last}; \mathbf{up})^+; \mathbf{right}. \quad (4)$$

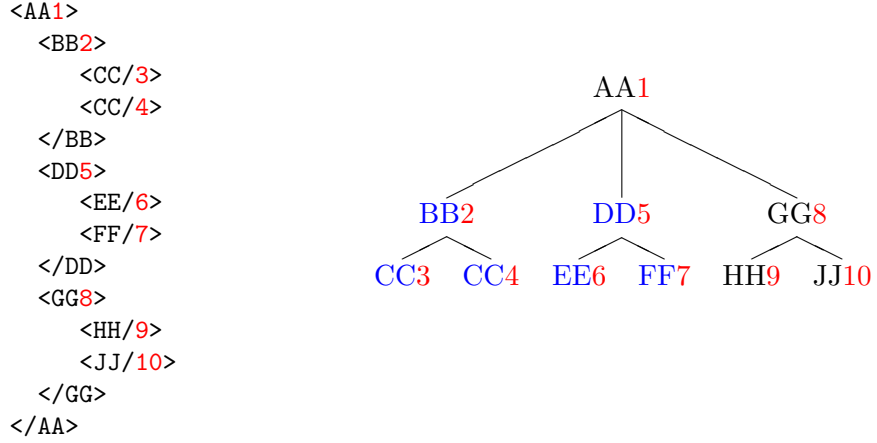


Figure 1: XML document and its DOM.

Now we show how to define the relation

$$x \ll y \wedge \phi(y) \wedge \forall z(x \ll z \ll y \rightarrow \phi(z)), \quad (5)$$

from the \mathcal{L}_{cp} programs. Having this it is easy to express the “until in document-order” modality: $\text{until}_{\ll}(\psi, \phi)$ holds at x iff $\exists yz((5) \wedge y \text{ down}; ?\psi z)$. Note that Theorem 4.1 ensures that this set is \mathcal{L}_{cp} definable, but not that the relation (5) is definable from the \mathcal{L}_{cp} programs.

We must use the definition of \mathcal{L}_{cp} programs containing union and composition. The definition is a case distinction based on the definition of $x \ll y$:

1. $x \text{ down}^+ y$
2. $x \text{ up}^+; \text{right}^+; \text{down}^+ y$
3. $x \text{ up}^+; \text{right}^+ y$
4. $x \text{ right}^+; \text{down}^+ y$
5. $x \text{ right}^+ y$.

We only show the easiest (first) and the hardest (second) case. The others are variations of these. For the first case we want to express that

$$x \text{ down}^+ y \wedge \phi(y) \wedge \forall z(x \ll z \ll y \rightarrow \phi(z)).$$

We explain our formulas by examples. Suppose x is node 1 and y is node 7 in Figure 1. Then ϕ must hold at nodes 2–7. This holds just in case

$$x (\text{down}; ?\phi \wedge [\text{left}^+][\text{down}^*]\phi)^+ y \quad (6)$$

holds.

For the second case we want to express that

$$x \text{ up}^+; \text{right}^+; \text{down}^+ y \wedge \phi(y) \wedge \forall z(x \ll z \ll y \rightarrow \phi(z)). \quad (7)$$

This holds exactly when x and y are related by

$$?[\text{down}^+] \phi \quad ; \quad (8)$$

$$(?[\text{right}^+][\text{down}^*] \phi; \text{up})^+ \quad ; \quad (9)$$

$$(\text{right}; ?[\text{down}^*] \phi)^* \quad ; \quad (10)$$

$$\text{right} ? \phi \quad ; \quad (11)$$

$$(\text{down}; ? \phi \wedge [\text{left}^+][\text{down}^*] \phi)^+ \quad (12)$$

This formula is best explained using a more elaborate tree, as in Figure 2. Suppose nodes C and R stand in the relation (7). Then (8) ensures that $\{A, B\}$ makes ϕ true; the test $[\text{right}^+][\text{down}^*] \phi$ in (9) will be evaluated at nodes C and G , thereby ensuring that ϕ holds in $\{F, D, E\}$ and $\{J, H, I\}$, respectively. The test $[\text{down}^*] \phi$ in (10) will be evaluated at all nodes strictly in between K and U , so here taking care that ϕ holds at $\{N, L, M\}$. (11) ensures that ϕ is true at U . Now (12) is just the formula (6) from the first case, ensuring that ϕ holds at $\{Q, O, P, T, R\}$.

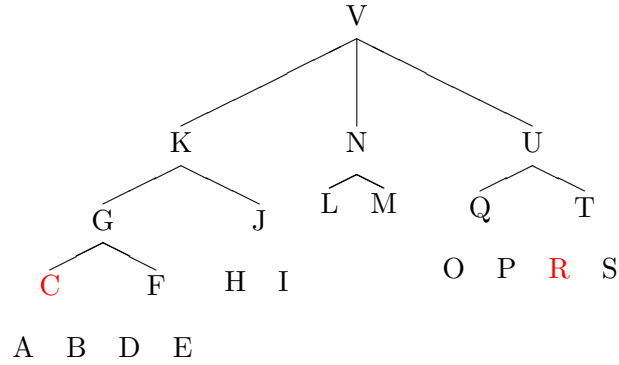


Figure 2: Example tree for the second case.

5 Expressivity 2: Second Order Properties

In this section we look at three concrete examples of non-trivial second order properties of trees that are expressible in PDL_{tree} ; first though, some background. The language PDL_{tree} can express properties beyond the reach of \mathcal{L}_{FO} . For example, PDL_{tree} can express the property of having an odd

number of daughters:

$$\langle \text{down} \rangle (\text{first} \wedge \langle (\text{right}; \text{right})^* \rangle \text{last}). \quad (13)$$

Note that the second conjunct $\langle (\text{right}; \text{right})^* \rangle \text{last}$ says that by chaining together a succession of double **right** steps we can reach the rightmost daughter node — which means that there must be an odd number of daughter nodes. This is not a property that any \mathcal{L}_{FO} formula can express.

On the other hand, PDL_{tree} is contained in $L_{K,P}^2$, Rogers monadic second-order logic of variably branching trees [27]. $L_{K,P}^2$ just extends \mathcal{L}_{FO} by quantification over unary predicates. The translation of PDL_{tree} formulas into $L_{K,P}^2$ is straightforward. Note, in particular, that we can use second-order quantification to define the transitive closure of a relation: for R any binary relation, xR^*y holds iff

$$x = y \vee \forall X (X(x) \wedge \forall z, z' (X(z) \wedge zRz' \rightarrow X(z')) \rightarrow X(y)).$$

This brings us to the most important open problem concerning PDL_{tree} :

Open Problem. Characterize the expressive power of PDL_{tree} interpreted on finite ordered trees in terms of a suitable fragment of monadic second order logic.

Within the context of query languages for XML documents a number of proposals for second order languages have been made. The goal then is to express unary MSO, MSO formulas denoting a set of nodes. We mention monadic datalog of [14] and the efficient tree logic of [22], which are both as expressive as unary MSO.

Neven and Schwentick [22] argue that unary MSO rather than \mathcal{L}_{FO} is the gold standard for a language designed for specifying nodes in finite ordered trees. Their most convincing example is a variant of the Boolean circuit problem. In order to obtain a better understanding of the second order expressivity of PDL_{tree} , we encode a number of second order properties in PDL_{tree} . In addition to the Boolean circuit problem just mentioned, we encode the frontier problem and we show that finiteness of ordered trees can be expressed in PDL_{tree} . The frontier problem is a typical linguistic problem. Expressing finiteness within a large class of tree like structures shows the robustness of the language. We look at the upshot of these examples at the end of this section. We start with the frontier problem.

5.1 The frontier problem

The frontier of a tree is the set of leaves ordered from left to right. In a parse tree of a natural language sentence, the frontier is exactly that sentence. Usually the frontier is where the actual data contained in a tree is located.

Given a condition ϕ on the frontier, we want to write an PDL_{tree} expression which is true at the root of a tree iff the frontier of the tree satisfies ϕ . For instance, ϕ could be a regular expression over variables, like $(p; q)^*$. The most natural application is when we know that each leaf node makes exactly one variable true. Then a tree satisfies ϕ if and only if the frontier is a word in $(p; q)^*$. But nothing forbids us to use arbitrary complex PDL_{tree} formulas in place of p and q . E.g., $\langle \text{up}^* \rangle \text{np}$ states that the first word of the parsed sentence (the first leaf of the frontier) is part of a noun-phrase (an **np**). Thus we do not view the frontier as a unique string, but as an infinite collection of strings, made up from formulas which are true at the respective nodes. Now let r be a regular expression in which arbitrary PDL_{tree} formulas are the letters. We say that a tree's frontier $l_1 \dots, l_n$ satisfies r iff there are PDL_{tree} formulas ϕ_i such that for all i , $l_i \models \phi_i$ and the string ϕ_1, \dots, ϕ_n is a word in r .

What we need for expressing the frontier condition is the successor relation between frontier nodes. This is naturally defined using the so called *document order* relation, abbreviated by \ll . For x, y nodes in a tree, we define

$$x \ll y \iff xR_{(\text{down}^+ \cup \text{up}^*; \text{right}^+; \text{down}^*)}y.$$

Figure 1 on page 9 exemplifies the document order, the numbers of the nodes correspond to their document ordering. The XML representation of the tree motivates the name of the ordering.

A frontier node y is the successor of a frontier node x if and only if $x \ll y$ and there is no leaf node in between x and y in the document order. An intuitive definition of the `next_frontier_node` relation between leaves can now be given as:

$$?\text{leaf}; \text{right}; ?\text{leaf} \cup ?\text{leaf}; (? \text{last}; \text{up})^+; \text{right}; (\text{down}; ? \text{first})^*; ?\text{leaf}. \quad (14)$$

Because we evaluate the PDL_{tree} formula at the root, we should add the step from the root to the first leaf to (14). So define the `next_frontier_node` relation as

$$?\text{root}; (\text{down}; ? \text{first})^*; ?\text{leaf} \cup (14).$$

Let `last_frontier_node` be an abbreviation of $\text{leaf} \wedge \langle (? \text{last}; \text{up})^* \rangle \text{root}$, which indeed is true exactly at the last frontier node (or simply at the root, if the model only has a root).

Now let r be a regular expression over a set of PDL_{tree} formulas. Then for any tree T , T 's frontier satisfies r if and only if the root of T satisfies $\langle r^\circ \rangle \text{last_frontier_node}$, where r° is r with $;$ placed between all PDL_{tree} formulas which act as letters in r and any such formula ϕ is replaced by

`next_frontier_node;?ϕ`. For instance, the frontier is in $(ab)^*$ iff the root satisfies

$$\langle (\text{next_frontier_node};?a;\text{next_frontier_node};?b)^* \rangle \text{last_frontier_node}.$$

Note that the formula is true on a tree containing only a root; thus it correctly recognizes the empty string.

5.2 The boolean circuit problem

We show how the boolean circuit problem can be expressed in PDL_{tree} . Our PDL_{tree} formula is based on the same idea as in [23]: use a depth first traversal of the tree. We start with defining the boolean circuit problem.

Definition 5.1 (Boolean circuits). Boolean circuits are finite $\{1, 0, C, D\}$ -labeled ordered binary trees such that

1. each leaf is labeled with exactly one of $\{1, 0\}$ and
2. each non-leaf is labeled with exactly one of $\{C, D\}$.

If \mathcal{B} is a boolean circuit and $b \in \mathcal{B}$ then with \mathcal{B}_b we denote the subtree of \mathcal{B} rooted at b . With $\mathcal{B}_{\bar{b}}$ we denote then tree which we obtain by removing everything below b . So in particular we have that b is a leaf of $\mathcal{B}_{\bar{b}}$.

The intended meaning of the labels is as one might expect: 1 means ‘true’, 0 means ‘false’, C means conjunction and D means disjunction. For any boolean circuit \mathcal{B} , define the boolean function `eval` from the domain of \mathcal{B} to $\{\text{‘true’}, \text{‘false’}\}$ in the expected way. For instance, as the Datalog program:

$$\begin{aligned} \text{eval}(x) & :- 1(x). \\ \text{eval}(x) & :- D(x), R_{\text{down}}(x,y), \text{eval}(y). \\ \text{eval}(x) & :- C(x), R_{\text{down}}(x,y), R_{\text{right}}(y,z), \text{eval}(y), \text{eval}(z). \end{aligned}$$

Also for any $b \in \mathcal{B}$ let $\text{height}(b)$ denote the length of the longest path starting at, but not including, b to a leaf. So if b is a leaf, then $\text{height}(b) = 0$.

General idea. To check if a boolean circuit evaluates to true we look for substructures that can be constructed as follows. We start at the root and move down. At disjunctive nodes we select one child. At conjunctive nodes we take both children. When we reach a leaf, it should be labeled with 1. We check if such a substructure exists in a depth first fashion. So, we walk down the tree, where at conjunctive nodes we always take the left route and make sure (by selecting the correct child at disjunctive nodes) we end up in a leaf labeled 1. We let the relation R_0 denote such a path. That is, for all x and y we have xR_0y iff the following three cases apply.

1. $\exists k \geq 1 t_1, \dots, t_k$ s.t. $x = t_1 \text{down} t_2 \text{down} \dots \text{down} t_k = y$

2. For all $1 \leq i < k$ if $t_i \Vdash C$, then $t_{i+1} \Vdash \text{first}$
3. $t_k \Vdash 1$

Next we walk up again until we are at a left child of a conjunctive node. We move right, to node b_r say, and repeat the procedure. When we return at node b_r we realize that we are about to enter a conjunctive node from the right and move further up until the next conjunctive node. With R_1 we denote this relation. So for all x and y we have xR_1y iff the following two cases apply.

1. $\exists k \geq 1 t_1, \dots, t_k$ s.t. $x = t_1 \text{up} t_2 \text{up} \dots \text{up} t_k = y$
2. For all $1 \leq i < k$, $t_i \Vdash \langle \text{up} \rangle C \rightarrow \text{last}$

When we reach the root of the boolean circuit the procedure stops. We can express both relations R_0 and R_1 by regular expressions π_0 and π_1 as follows. Let π_0 be the regular expression which corresponds to R_0 . That is

$$\pi_0 = ((?D; \text{down}) \cup (?C; \text{down}; ?\text{first}))^*; ?1.$$

Let π_1 be the regular expression corresponding to R_1 . That is

$$\pi_1 = (?(\langle \text{up} \rangle C \rightarrow \text{last}) \wedge \neg \text{root}; \text{up})^*.$$

Finally define

$$\beta = \langle \pi_0; \pi_1; (\text{right}; \pi_0; \pi_1)^* \rangle \text{root}.$$

Before we move on let us make a remark. On first sight one might think that we need in the definition of R_1 a third clause. Namely

3. $t_k \Vdash \langle \text{up} \rangle C \wedge \neg \text{last}$ or $t_k \Vdash \text{root}$.

And, consequently, in stead of π_1 we should have

$$\pi_1; ?(\langle \text{up} \rangle C \wedge \neg \text{last}) \vee \text{root}.$$

This is not necessary. With the current definition of R_1 we allow for a check (but do not consider it necessary) that the second child of a disjunctive node is true when we already know that the first child is. This is just as harmless as it is useless. Nevertheless the proof below (in particular Lemma 5.5) does not work without this omission.

Theorem 5.2. *β is forced at the root r of a boolean circuit iff. $\text{eval}(r)$ is true.*

Lemma 5.3. *Let \mathcal{B} be a boolean circuit. For all nodes $b \in \mathcal{B}$ we have the following.*

1. $b \Vdash \beta \wedge C \rightarrow [\mathbf{down}]\beta$

2. $b \Vdash \beta \wedge D \rightarrow \langle \mathbf{down} \rangle \beta$

Proof. First we show 1. Suppose $b \Vdash \beta \wedge C$. Let b_l be the left child of b and b_r be the right child of b . It is easy to see that $b_l \Vdash \beta$. To show that $b_r \Vdash \beta$ we need a lemma.

Lemma. *For any x for which not $x(\mathbf{up})^*b_l$ (e.a. $x \notin \mathcal{B}_{b_l}$) we have that if $b_l(\pi_0; \pi_1; (\mathbf{right}; \pi_0; \pi_1)^*x)$ then $b_l(\mathbf{right}; \pi_0; \pi_1)^*x$.*

Proof. Choose x as stated. We show with induction on n that

if $b_l(\pi_0; \pi_1; (\mathbf{right}; \pi_0; \pi_1)^n x)$ then $b_l(\mathbf{right}; \pi_0; \pi_1)^*x$.

If $n = 0$ then for some t , $b_l\pi_0t\pi_1x$. Clearly $t(\mathbf{up})^*b_l$ and $t(\mathbf{up})^*x$. So, by choice of x , $b_l(\mathbf{up})^+x$. But this is clearly in contradiction with the definition of π_1 .

Now suppose $b_l(\pi_0; \pi_1; (\mathbf{right}; \pi_0; \pi_1)^{n+1}x)$. Choose t such that

$b_l(\pi_0; \pi_1; (\mathbf{right}; \pi_0; \pi_1)^n t)$ and $t(\mathbf{right}; \pi_0; \pi_1)x$.

We can assume that $t(\mathbf{up})^*b_l$ (otherwise we are done by (IH)). We also can assume that $t \neq b_l$ and thus $t(\mathbf{up})^+b_l$. Fix some t' for which $t(\mathbf{right}; \pi_0)t'\pi_1x$. By the above we obtain $t'(\mathbf{up})^+b_l$. Similar as in the case $n = 0$ this leads us to a contradiction. ⊥

Now we continue with showing that $b_r \Vdash \beta$. Since $b_l \Vdash \beta$ we can find some x_1, x_2, \dots such that

$b_l = x_1(\pi_0; \pi_1)x_2(\mathbf{right}; \pi_0; \pi_1)x_3 \cdots (\mathbf{right}; \pi_0; \pi_1)r$.

Where r is the root of \mathcal{B} . Let i be that smallest number such that not $x_i(\mathbf{up})^*b_l$. Note that $i > 2$. So, by the above lemma and by choice of i , we have $b_l(\mathbf{right}; \pi_0; \pi_1)x_i$. So, $b_r(\pi_0; \pi_1)x_i$ and thus $b_r \Vdash \beta$. We have shown 1.

Item 2 is rather trivial. For if we suppose that $b \Vdash \beta \wedge D$ then it is easy to verify, using the definition of π_0 , that $b \Vdash \langle \mathbf{down} \rangle \beta$. ⊥

Corollary 5.4. *Let \mathcal{B} be a boolean circuit. For all nodes $b \in \mathcal{B}$ we have that if $b \Vdash \beta$ then $\text{eval}(b)$ is true.*

Proof. Induction on $\text{height}(b)$. If $\text{height}(b) = 0$ then the claim is clear by the definition of π_0 . So suppose $\text{height}(b) > 0$. There are two cases to consider.

Case: $b \Vdash C$. By Lemma 5.3 we have $b \Vdash [\mathbf{down}]\beta$. So, by (IH), we have that for all children b' of b that $\text{eval}(b')$ is true and thus $\text{eval}(b)$ is true.

Case: $b \Vdash D$. By Lemma 5.3 we have $b \Vdash \langle \text{down} \rangle \beta$. So, by (IH), for some child b' of b we have $\text{eval}(b')$ is true and thus $\text{eval}(b)$ is true. \dashv

Lemma 5.5. *Let \mathcal{B} be a boolean circuit. For all $b \in \mathcal{B}$ for which $\text{eval}(b)$ is true we have that $b(\pi_0; \pi_1; (\mathbf{right}; \pi_0; \pi_1)^*)b$.*

Proof. induction on $\text{height}(b)$. If $\text{height}(b) = 0$ then this is clear. So suppose $\text{height}(b) > 0$. There are two cases to consider.

Case: $b \Vdash C$. Then for both children b_l and b_r of b we have $\text{eval}(b_l)$ is true and $\text{eval}(b_r)$ is true. By (IH), $b_l(\pi_0; \pi_1; (\mathbf{right}; \pi_0; \pi_1)^*)b_l$ and $b_r(\pi_0; \pi_1; (\mathbf{right}; \pi_0; \pi_1)^*)b_r$. So,

$$b(?C; \text{down}; ?\text{first}; \pi_0; \pi_1; (\mathbf{right}; \pi_0; \pi_1)^*; \\ \mathbf{right}; \pi_0; \pi_1; (\mathbf{right}; \pi_0; \pi_1)^*; ?\text{last}; \text{up})b.$$

Thus, as one can easily verify, we have

$$b(\pi_0; \pi_1; (\mathbf{right}; \pi_0; \pi_1)^*)b.$$

Case: $b \Vdash D$. Then for at least one child b_i of b we have $\text{eval}(b_i)$ is true. So by (IH) we obtain $b_i(\pi_0; \pi_1; (\mathbf{right}; \pi_0; \pi_1)^*)b_i$. Thus

$$b(?D; \text{down}; (\pi_0; \pi_1; (\mathbf{right}; \pi_0; \pi_1)^*); ?\neg\langle \text{up} \rangle C; \text{up})b.$$

Which implies $b(\pi_0; \pi_1; (\mathbf{right}; \pi_0; \pi_1)^*)b$. \dashv

Now we are ready to prove Theorem 5.2.

Proof of Theorem 5.2. (\Rightarrow) Immediate from Corollary 5.4. (\Leftarrow) Suppose $\text{eval}(r)$ is true. By Lemma 5.5 we have $r(\pi_0; \pi_1; (\mathbf{right}; \pi_0; \pi_1)^*)r$. So in particular $r \Vdash \langle \pi_0; \pi_1; (\mathbf{right}; \pi_0; \pi_1)^* \rangle \text{root}$, e.a. $r \Vdash \beta$. \dashv

5.3 Expressing finiteness

In this subsection we let go the restriction to finite trees. Normally one would define arbitrary trees as partially ordered sets $(W, <)$ with a unique root and such that for each $w \in W$ the set $\{v \mid v < w\}$ is well-ordered by $<$. The height of a node w is then defined as the ordertype of $\{v \mid v < w\}$ and we say that a tree is of height ω when the height of each node is finite. We can do a little bit better. Below we define first order definable structures such that the part that PDL_{tree} can see is a tree of height ω .

Fist, for a binary relation R we say that y is a *direct successor* of x when xRy and for no z we have $xRzRy$. We define *direct predecessor* in a similar way. We say that R is *discrete* when for any xRy such that y is not a direct successor of x , there exists some direct successor z of x with $xRzRy$. Notice that discrete relations are always irreflexive. We say that a structure $\langle T, R_{\text{down}^+}, R_{\text{right}^+} \rangle$ (Note that in this context, R_{down^+} and R_{right^+} are primitive relation symbols themselves.) is *tree-like* when

1. R_{down^+} is a discrete and partial order on T with a unique root,
2. each $t \in T$ has at most one direct R_{down^+} -predecessor,
3. R_{right^+} is discrete and linearly orders the direct successors of any $t \in T$, in particular if $xR_{\text{right}^+}y$ then x and y have the same direct R_{down^+} -predecessor.

Clearly this class of structures is first order definable within the class of Kripke frames with two accessibility relations and any tree is a tree-like structure. We define the relations R_{down} , R_{right} and all the other relations R_π that may occur within PDL_{tree} -modalities as above in Section 2. But please note that although we do have $R_{\text{down}}^+ \subseteq R_{\text{down}^+}$, $R_{\text{right}}^+ \subseteq R_{\text{right}^+}$, in general these inclusions will be proper. If $\mathcal{T} = \langle T, R_{\text{down}^+}, R_{\text{right}^+} \rangle$ is a tree-like structure with root r , then we write \mathcal{T}_r for the structure $\langle r \rangle_{\langle T, R_{\text{down}}, R_{\text{right}} \rangle}$, the substructure of $\langle T, R_{\text{down}}^+, R_{\text{right}}^+ \rangle$ generated by r using the defined relations R_{down} and R_{right} , in the usual modal logic sense. Of course for any PDL_{tree} formula ϕ we have that $\mathcal{T}, r \Vdash \phi$ iff. $\mathcal{T}_r, r \Vdash \phi$. So without danger of confusion we can write $r \Vdash \phi$.

As a sort of corollary to the proof of the definability of boolean circuits we will show here that PDL_{tree} can define finiteness of tree-like structures.

Theorem 5.6. *There exists a PDL_{tree} -formula Fin such that for any tree-like structure \mathcal{T} with root r we have $\mathcal{T}, r \Vdash \text{Fin}$ iff. \mathcal{T} is finite.*

Proof. Let δ and γ be as defined in (15) and (16) below and let Fin be $\delta \wedge \gamma$. The proof proceeds in stages. In Lemma 5.7 we show that it is sufficient to show that $\mathcal{T}_r, r \Vdash \text{Fin}$ iff. \mathcal{T}_r is finite. This latter task is performed in the Lemmata 5.8, 5.9 and 5.10 ◻

Lemma 5.7. *For any tree-like structure \mathcal{T} with root r , \mathcal{T}_r is an ordered² tree of height ω , and \mathcal{T}_r is finite iff. \mathcal{T} is finite.*

Proof. The first assertion is a direct consequence of the definition of tree-like structures. The second assertion follows from the fact that if x is a leaf in \mathcal{T}_r then by discreteness there does not exist any R_{down^+} descendant of x in \mathcal{T} . ◻

As a first approximation for finiteness put

$$\delta = [\text{down}^*](\langle \text{left}^* \rangle \text{first} \wedge \langle \text{right}^* \rangle \text{last}). \quad (15)$$

Lemma 5.8. *For any tree-like structure \mathcal{T} with root r we have that $\mathcal{T}_r, r \Vdash \delta$ iff. \mathcal{T}_r is finitely branching.*

²In case the tree is infinitely branching the sibling order R_{right} might be non-total, but no matter, see Lemma 5.8.

Proof. The left to right direction holds since if $t \in \mathcal{T}_r$ has infinitely many children then by discreteness we can find an infinite, to the left or to the right, R_{right} -chain. The converse is obvious. \dashv

So in order to define the class of finite tree-like structures it is enough to define the class of finite trees as a subclass from the class of ordered trees of height ω which are finitely branching. To this end put

$$\begin{aligned}\pi_0 &= (\text{down}; ?\text{first})^*; ?\text{leaf}, \\ \pi_1 &= (?!\text{last}; \text{up})^*\end{aligned}$$

and

$$\gamma = \langle \pi_0; \pi_1; (\text{right}; \pi_0; \pi_1)^* \rangle \text{root}. \quad (16)$$

Before we move on let us introduce some terminology. A *branch* b in a finitely branching tree T of height ω is a sequence

$$r = x_1(\text{down})x_2(\text{down}) \cdots (\text{down})x_n(\text{down}) \cdots$$

where r is the root of T and either b is infinite and in this case is (down) closed, or its last element is a leaf. If b and b' are branches then we say that b is *to the left* of b' whenever if i is the smallest i such that $b_i \neq b'_i$ then $b_i(\text{right})^+ b'_i$. Clearly, since T is finitely branching, this gives us a linear ordering on the branches of some fixed tree. For $t \in T$ and b a branch of T we write $t < b$ if $t \notin b$ and t occurs on some branch to the left of b . $t \leq b$ means $t < b$ or $t \in b$. Similar definitions hold for $b < t$, $b \leq t$.

Lemma 5.9. *Suppose T is finitely branching tree of height ω and $t \in T$. If T_t is finite then $t(\pi_0; \pi_1; (\text{right}; \pi_0; \pi_1)^*)t$.*

Proof. Induction on $\text{height}(t)$. Similar to the proof of Lemma 5.5. \dashv

Lemma 5.10. *Suppose T is finitely branching tree of height ω with root r . If T is infinite then not $r(\pi_0; \pi_1; (\text{right}; \pi_0; \pi_1)^*)r$.*

Proof. Since T is infinite, finitely branching and of height ω , T must contain an infinite branch. Let b be the leftmost infinite branch of T . Such a branch can easily be constructed by starting from r and in each successive step select the leftmost child of the previously selected node which roots an infinite subtree. The following is obvious.

1. $x \leq b$ and $x(\pi_0)y$ imply $y < b$
2. $x < b$ and $x(\pi_1)y$ imply $y < b$
3. $x < b$ and $x(\text{right})y$ imply $y \leq b$

1 is clear, since π_0 only walks to leftmost children. 2 is clear, since π_1 only walks from rightmost children. 3 is clear, by definition of the ordering on branches.

Now let us assume that $r(\pi_0; \pi_1; (\mathbf{right}; \pi_0; \pi_1)^*)r$. Then there exists some sequence

$$r = a_0(\pi_0)a_1(\pi_1)a_2(\mathbf{right})a_3(\pi_0)a_4 \cdots a_{k-2}(\pi_1)a_{k-1} = r.$$

By the above three points it follows, with induction on i , that

$$\text{for all } i < k, a_i \leq b. \tag{17}$$

Since $a_{k-3}(\pi_0)a_{k-2}$ we have that a_{k-2} must be a leaf, and since $a_{k-2}(\pi_1)t$ we also have that the branch in T ending in a_{k-2} only contains rightmost nodes. But this implies that b , as the leftmost infinite branch of T , must be on the left of the branch ending in a_{k-2} . So in particular $b \leq a_{k-2}$. But since a_{k-2} is a leaf we even have $b < a_{k-2}$, in contradiction with (17). \dashv

5.4 The upshot

What is the upshot of these examples? First and foremost they were intended to show the language in action. To show that semantic reasoning is naturally captured in PDL_{tree} formulas, even when it comes to hard problems. Even though we provided rigorous correctness proofs, we feel that once the semantic argument is understood, correctness of the PDL_{tree} formalisation is almost self evident.

Although Boolean circuits looks like a canonical MSO problem it has certain peculiarities which we could exploit, in particular that one depth-first traversal of the tree is sufficient to determine the truth of the formula. The problem suggest a possible strengthening of the language: intersection of programs with $?\top$. With this we can specify the set of all points t at which $\text{eval}(t)$ is true, and not just the root.

6 Conclusions

We hope that we convinced the reader that PDL is a very natural formalism for reasoning about ordered trees. Of course we could do all we did in monadic second order logic, but the absence of variables, and the restriction of the program connectives to the regular expression operators has a number of advantages. We mention three, which are all speculative and debatable.

Firstly, we believe that using intersection and complementation in creating relations is difficult, and gives rise to specifications whose down-to-earth meaning is not immediately obvious. The regular expression operators stay much closer to the semantics.

Secondly, there is the difference in complexity. Consider the sublanguage \mathcal{L}_{cp} , which is equally expressive as first order logic \mathcal{L}_{FO} on ordered trees. The satisfiability problem for \mathcal{L}_{FO} (and a fortiori for MSO) is hard for non elementary time, while the equally expressive \mathcal{L}_{cp} can be decided in single exponential time. This is an enormous improvement of course, and makes one think that \mathcal{L}_{cp} is much better. On the other hand, \mathcal{L}_{FO} is non elementary more succinct than \mathcal{L}_{cp} , meaning that there are \mathcal{L}_{FO} formulas whose smallest equivalent \mathcal{L}_{cp} formula has size roughly a tower of 2's whose length is bounded by the quantifier depth of the first order formula. From this perspective one could prefer \mathcal{L}_{FO} as a more user or programmer friendly language. On the other hand, writing \mathcal{L}_{cp} formulas seems more honest: compare the first order characterization of until in document order versus the \mathcal{L}_{cp} characterization. The latter is close to programming a tree automaton, and gives a plan how to check whether the formula is true on a specific tree. The first order formula does not provide a clue how to evaluate it.

We also strongly believe that languages without variables are easier to work with. Independent evidence for this comes from the W3C endorsed language XPath 1.0 [7] whose navigational version [15] is almost as expressive as \mathcal{L}_{Core} [21] and has virtually the same syntax, also without variables. XPath 1.0 is language designed for selecting nodes from XML documents. XPath plays a crucial role in other XML technologies such as XSLT [10], XQuery [9] and XML schema constraints, e.g., [8].

Let us return to the main open problem of the paper. Several proposals for unary MSO complete languages have been made, but none of these is as simple and easy to learn as PDL_{tree} . So the most pressing open problem seems to be to determine the exact expressive power of PDL_{tree} compared to unary MSO and —assuming there is a difference— to determine whether the extra expressivity given by unary MSO is useful in specific applications as linguistics or the XML-world.

Acknowledgments. Thanks are due to Michael Benedikt for suggesting the frontier and the boolean circuit problem, and to Wouter Kuijpers for the formula expressing finiteness, and to Jan Hidders and Frank Neven for valuable feedback. Part of this work was carried out as part of the INRIA funded research partnership between LIT (Language and Inference Technology Group, University of Amsterdam) and LED (Langue et Dialogue, LORIA, Nancy). Marx and Goris are supported by NWO grant 612.000.106 and Afanasiev by NWO grant 612.000.207. De Rijke is supported by grants from NWO, under project numbers 365-20-005, 612.069.006, 612.000.106, 220-80-001, 612.000.207, and 612.066.302.

References

- [1] S. Abiteboul, P. Buneman, and D. Suciu. *Data on the web*. Morgan Kaufman, 2000.
- [2] P. Blackburn, M. de Rijke, and Y. Venema. *Modal Logic*. Cambridge University Press, 2001.
- [3] P. Blackburn, B. Gaiffe, and M. Marx. Variable free reasoning on finite trees. In *Proceedings of Mathematics of Language (MOL-8)*, Bloomington, 2003.
- [4] P. Blackburn and W. Meyer-Viol. Linguistics, logic, and finite trees. *Logic Journal of the IGPL*, 2:3–29, 1994.
- [5] P. Blackburn, W. Meyer-Viol, and M. de Rijke. A proof system for finite trees. In H. Kleine Büning, editor, *Computer Science Logic*, volume 1092 of *LNCS*, pages 86–105. Springer, 1996.
- [6] World-Wide Web Consortium. Extensible Markup Language (XML) 1.0 (second edition) W3C. <http://www.w3.org/TR/REC-xml>.
- [7] World-Wide Web Consortium. XML path language (XPath): Version 1.0. <http://www.w3.org/TR/xpath.html>.
- [8] World-Wide Web Consortium. XML schema part 1: Structures. <http://www.w3.org/TR/xmlschema-1>.
- [9] World-Wide Web Consortium. XQuery 1.0: A query language for XML. <http://www.w3.org/TR/xquery/>.
- [10] World-Wide Web Consortium. XSL transformations language (XSLT): Version 2.0. <http://www.w3.org/TR/xslt20/>.
- [11] K. Etessami, M. Vardi, and Th. Wilke. First-order logic with two variables and unary temporal logic. In *Proceedings 12th Annual IEEE Symposium on Logic in Computer Science*, pages 228–235, Warsaw, Poland, 1997. IEEE.
- [12] M. Fisher and R. Ladner. Propositional dynamic logic of regular programs. *J. Comput. Syst. Sci.*, 18(2):194–211, 1979.
- [13] D.M. Gabbay and F. Guenther, editors. *Handbook of Philosophical Logic*, volume 2. Reidel, Dordrecht, 1984.
- [14] G. Gottlob and C. Koch. Monadic datalog and the expressive power of languages for web information extraction. In *Symposium on Principles of Database Systems (PODS)*, pages 17–28, 2002.
- [15] G. Gottlob, C. Koch, and R. Pichler. Efficient algorithms for processing XPath queries. In *Proc. of the 28th International Conference on Very Large Data Bases (VLDB 2002)*, 2002.
- [16] D. Harel, D. Kozen, and J. Tiuryn. *Dynamic Logic*. MIT Press, 2000.
- [17] J.A.W. Kamp. *Tense Logic and the Theory of Linear Order*. PhD thesis, University of California, Los Angeles, 1968.
- [18] M. Kracht. Syntactic codes and grammar refinement. *Journal of Logic, Language and Information*, 4:41–60, 1995.

- [19] M. Kracht. Inessential features. In Christian Retore, editor, *Logical Aspects of Computational Linguistics*, number 1328 in LNAI, pages 43–62. Springer, 1997.
- [20] M. Marx. Conditional XPath, the first order complete XPath dialect. In *Proceedings of PODS'04*, 2004.
- [21] M. Marx and M. de Rijke. Semantic characterizations of XPath. In *Submitted manuscript*, 2004.
- [22] F. Neven and T. Schwentick. Expressive and efficient pattern languages for tree-structured data. In *Proc. PODS*, pages 145–156. ACM, 2000.
- [23] F. Neven and Th. Schwentick. On the power of tree-walking automata. In E. Wel U. Montanari, J. Rolim, editor, *Automata, Languages and Programming, Proceedings, ICALP 2000, LNCS 1853*, pages 547–560. Springer, 2000.
- [24] A. Palm. *Transforming tree constraints into formal grammars*. PhD thesis, Universität Passau, 1997.
- [25] A. Palm. Propositional tense logic for trees. In *Sixth Meeting on Mathematics of Language*. University of Central Florida, Orlando, Florida, 1999.
- [26] V. Pratt. Models of program logics. In *Proceedings of the 20th IEEE symposium on Foundations of Computer Science*, pages 115–122, 1979.
- [27] J. Rogers. *A Descriptive Approach to Language Theoretic Complexity*. CSLI Press, 1998.
- [28] E. Spaan. *Complexity of modal logics*. PhD thesis, University of Amsterdam, Institute for Logic, Language and Computation, 1993.
- [29] M.Y. Vardi and P. Wolper. Automata-theoretic techniques for modal logics of programs. *Journal of Computer and System Sciences*, 32:183–221, 1986.
- [30] M. Weyer. Decidability of S1S and S2S. In E. Grädel et al., editor, *Automata, Logics, and Infinite Games*, volume 2500 of *LNCS*, pages 207–230. Springer, 2002.