

First Steps Towards Taming Description Logics with Strings

Stéphane Demri

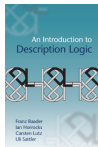
Joint work with Karin Quaas
(Leipzig University, Germany)

GT FM & AI, January 2024

Reasoning on Ontologies with Description Logics

- Ontology: formal specification of some domain with concepts, objects, relationships between concepts, objects, etc.
- Backbone of ontologies includes:
 - taxonomy (classification of objects),
 - axioms (to constrain the models of the defined terms).
- Description logics are well-known logical formalisms dedicated to ontologies.

[Baader et al, Book 2017]



- BioPortal (<http://bioportal.bioontology.org/>): huge amount of ontologies such as Cell Ontology and Plant Ontology to facilitate scientific activities.
- Computational properties.
 - Acceptable trade-off between expressivity and complexity.
 - Decidability and often tractability.
 - Implementation in tools of the main reasoning tasks.

Description Logics with Concrete Domains

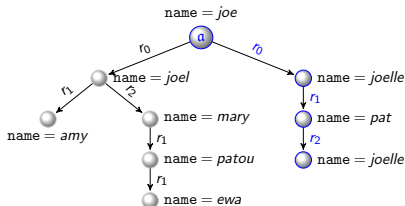
- Need to express concrete properties about data in ontologies. (e.g. age, duration, name, size, etc.)
- Concrete domain $\mathcal{D} = (\mathbb{D}, R_1, R_2, \dots)$: fixed non-empty domain with a family of relations.
- $(\mathbb{N}, <, +1)$, $(\mathbb{Q}, <, =)$, $(\mathbb{N}, <, =)$, $(\{0, 1\}^*, <_{\text{pre}}, <_{\text{suf}})$.
- Concrete domain RCC8 with space regions in \mathbb{R}^2 contains topological relations between spatial regions.
See e.g. [Wolter & Zakharyashev, KR'00]
- General scheme for integrating concrete domains in DLs.
[Baader & Hanschke, IJCAI'91]
 - declarative semantics close to the usual semantics for DLs,
 - generic extensions of DLs with various concrete domains,
 - tableaux-based algorithms combined with theory reasoning.

Methods for Handling Concrete Domains

- Tableaux-based decision procedures for ω -admissible concrete domains. [Lutz & Miličić, JAR 2007]
 - $\mathcal{R} = (\mathbb{R}, <, =, (=_r)_{r \in \mathbb{R}})$ is ω -admissible.
 - $\mathcal{N} = (\mathbb{N}, <, =, (=_n)_{n \in \mathbb{N}})$ is not ω -admissible.
- Translation into a decidable extension of MSO with bounding quantifier B. [Carapelle & Turhan, ECAI'16]
 - EHD approach developed with $\text{Bool}(\text{MSO}, \text{WMSO} + \text{B})$ over infinite trees of finite branching degree. [Carapelle & Kartzow & Lorhey, JCSS 2016]
 - Decidability of concept satisfiability problem w.r.t. general TBoxes for $\mathcal{ALC}(\mathcal{N})$.
- Translation into Rabin tree automata over finite alphabets using approximations for satisfiable symbolic interpretations.
 - Concept satisfiability problem w.r.t. general TBoxes for $\mathcal{ALC}(\mathcal{N})$ in EXPTIME. [Labai & Ortiz & Šimkus, KR'20]

Finite Strings with the Prefix Relation

- $\mathcal{D}_\Sigma = (\Sigma^*, <_{\text{pre}}, =, (=w)_{w \in \Sigma^*})$.
- $(\mathbb{N}, <, =, (=n)_{n \in \mathbb{N}})$ corresponds to \mathcal{D}_Σ with singleton Σ .



$$(\exists r_0 r_1 r_2 \cdot (\text{name} <_{\text{pre}} SSS \text{name})) \wedge \forall r_0 \cdot (\text{name} <_{\text{pre}} S \text{name})$$

('S' similar to next in temporal logic)

- Concept satisfiability problem w.r.t. general TBoxes for $\mathcal{ALC}(\mathcal{W})$ with $\mathcal{W} = (\Sigma^*, \cdot, =, (=w)_{w \in \Sigma^*})$ is undecidable.

[Lutz, PhD 2002]

- No known decidability results for description logics with \mathcal{D}_Σ .

ALC in a Nutshell

- Complex concepts.

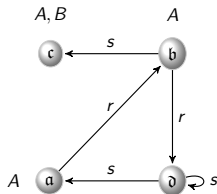
$$C ::= \top \mid \perp \mid A \mid \neg C \mid C \sqcap C \mid C \sqcup C \mid \exists r.C \mid \forall r.C,$$

with concept names A and role names r .

- Interpretation $\mathcal{I} \stackrel{\text{def}}{=} (\Delta^{\mathcal{I}}, \cdot^{\mathcal{I}})$

- $\Delta^{\mathcal{I}}$: non-empty set (the *domain*).
- $\cdot^{\mathcal{I}}$: *interpretation function* such that

$$A^{\mathcal{I}} \subseteq \Delta^{\mathcal{I}} \quad r^{\mathcal{I}} \subseteq \Delta^{\mathcal{I}} \times \Delta^{\mathcal{I}}$$



Concept name A /role name r

\approx

unary predicate/binary predicate

Set-Theoretical Semantics for Complex Concepts

$$\top^{\mathcal{I}} \stackrel{\text{def}}{=} \Delta^{\mathcal{I}}$$

$$\perp^{\mathcal{I}} \stackrel{\text{def}}{=} \emptyset$$

$$(\neg C)^{\mathcal{I}} \stackrel{\text{def}}{=} \Delta^{\mathcal{I}} \setminus C^{\mathcal{I}}$$

$$(C_1 \sqcup C_2)^{\mathcal{I}} \stackrel{\text{def}}{=} C_1^{\mathcal{I}} \cup C_2^{\mathcal{I}}$$

$$(C_1 \sqcap C_2)^{\mathcal{I}} \stackrel{\text{def}}{=} C_1^{\mathcal{I}} \cap C_2^{\mathcal{I}}$$

$$(\exists r.C)^{\mathcal{I}} \stackrel{\text{def}}{=} \{a \in \Delta^{\mathcal{I}} \mid r^{\mathcal{I}}(a) \cap C^{\mathcal{I}} \neq \emptyset\}$$

$$(\forall r.C)^{\mathcal{I}} \stackrel{\text{def}}{=} \{a \in \Delta^{\mathcal{I}} \mid r^{\mathcal{I}}(a) \subseteq C^{\mathcal{I}}\}$$

$$\mathcal{R}(a) \stackrel{\text{def}}{=} \{b \mid (a, b) \in \mathcal{R}\}$$

Inclusion and Decision Problem TSAT(\mathcal{ALC})

- General concept inclusions $C \sqsubseteq D$ (GCIs).
E.g., $\text{Employee} \sqsubseteq \exists \text{WorksFor}.$

$$\mathcal{I} \models C \sqsubseteq D \stackrel{\text{def}}{\Leftrightarrow} C^{\mathcal{I}} \subseteq D^{\mathcal{I}}$$

- Terminological Box (TBox) \mathcal{T} : finite collection of GCIs.
- Interpretation $\mathcal{I} = (\Delta^{\mathcal{I}}, \cdot^{\mathcal{I}})$, TBox \mathcal{T} .

$$\mathcal{I} \models \mathcal{T} \stackrel{\text{def}}{\Leftrightarrow} \text{for all } C \sqsubseteq D \in \mathcal{T}, \mathcal{I} \models C \sqsubseteq D$$

- *Concept satisfiability problem w.r.t. general TBoxes* (TSAT(\mathcal{ALC})):

Input: A concept C_0 and a TBox \mathcal{T} .

Question: Is there an interpretation \mathcal{I} such that $\mathcal{I} \models \mathcal{T}$
and $C_0^{\mathcal{I}} \neq \emptyset$?

- TSAT(\mathcal{ALC}) is EXPTIME-complete.

Description Logic $\mathcal{ALCF}^{\mathcal{P}}(\mathcal{D}_{\Sigma})$

$$\exists r_0 r_1 r_2 \cdot (\text{name} <_{\text{pre}} \text{SSS name})$$

- New (atomic) concepts of the form $\exists P. [\Theta]$ and $\forall P. [\Theta]$:
 - non-empty sequence P of role names (*role path*),
 - Boolean constraint Θ built over terms of the form $S^j \mathbf{x}$ with $j \leq |P|$ and atomic constraints of the form

$$t <_{\text{pre}} t' \quad t = t' \quad =_{\mathfrak{w}}(t) \text{ (also written } t = \mathfrak{w})$$

- Interpretation $\mathcal{I} = (\Delta^{\mathcal{I}}, \cdot^{\mathcal{I}}, \mathfrak{v})$ with $\mathfrak{v}: \Delta^{\mathcal{I}} \times \text{VAR} \rightarrow \Sigma^*$.
- $(r_1 r_2 \dots r_n)^{\mathcal{I}} \stackrel{\text{def}}{=} \text{set of tuples } (\mathfrak{a}_0, \dots, \mathfrak{a}_n) \text{ in } (\Delta^{\mathcal{I}})^{n+1} \text{ such that } (\mathfrak{a}_{i-1}, \mathfrak{a}_i) \in r_i^{\mathcal{I}} \text{ for all } i \in [1, n].$
- Satisfaction relation $\mathcal{I}, \pi = (\mathfrak{a}_0, \mathfrak{a}_1, \dots, \mathfrak{a}_n) \models \Theta$:



- $\mathcal{I}, \pi \models S^i \mathbf{x} <_{\text{pre}} S^j \mathbf{y} \stackrel{\text{def}}{\Leftrightarrow} \mathfrak{v}(\mathfrak{a}_i, \mathbf{x}) <_{\text{pre}} \mathfrak{v}(\mathfrak{a}_j, \mathbf{y}),$ (similar for $=$)
- $\mathcal{I}, \pi \models S^i \mathbf{x} =_{\mathfrak{w}} \stackrel{\text{def}}{\Leftrightarrow} \mathfrak{v}(\mathfrak{a}_i, \mathbf{x}) =_{\mathfrak{w}},$

Clauses for Interpreting the Concepts Involving \mathcal{D}_Σ

$$(\exists P. [\Theta])^{\mathcal{I}} \stackrel{\text{def}}{=} \{a_0 \in \Delta^{\mathcal{I}} \mid \exists a_1, \dots, a_n \in \Delta^{\mathcal{I}} \text{ s.t. } \pi = (a_0, a_1, \dots, a_n) \in P^{\mathcal{I}} \text{ and } \mathcal{I}, \pi \models \Theta\}$$

$$(\forall P. [\Theta])^{\mathcal{I}} \stackrel{\text{def}}{=} \{a_0 \in \Delta^{\mathcal{I}} \mid \forall a_1, \dots, a_n \in \Delta^{\mathcal{I}}, \pi = (a_0, a_1, \dots, a_n) \in P^{\mathcal{I}} \text{ implies } \mathcal{I}, \pi \models \Theta\}$$

- $\mathcal{ALCF}^{\mathcal{P}}(\mathcal{D}_\Sigma)$ has also functional role names, omitted today.
- Integrating concrete domains may come with some extra costs: we need to solve a (potentially infinite) constraint satisfaction problem.
- No finite interpretation property, e.g. with

$$\mathcal{I} = \{\top \sqsubseteq \exists r. [\mathbf{x} <_{\text{pre}} S\mathbf{x}]\}$$

Main Steps For Getting ExpTime Upper Bound

- Automata-based approach with tree constraint automata (TCA) accepting infinite data trees with domain Σ^* .
- Step 0: to transform the instance so that every concept is in *simple form*, proper form to perform Step 1.
- Step 1: $C_0, \mathcal{T} \longrightarrow \mathbb{A}$ s.t. C_0, \mathcal{T} is a positive instance of TSAT iff $L(\mathbb{A}) \neq \emptyset$.

- Step 2: \mathbb{A} (for \mathcal{D}_Σ) $\longrightarrow \mathbb{A}'$ (for \mathcal{N}) such that

$$L(\mathbb{A}) \neq \emptyset \text{ iff } L(\mathbb{A}') \neq \emptyset$$

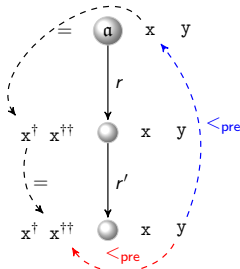
- Step 3: complexity analysis to get EXPTIME-completeness of TSAT. (based on [Demri & Quaas, CONCUR'23]).

Concepts in Simple Form

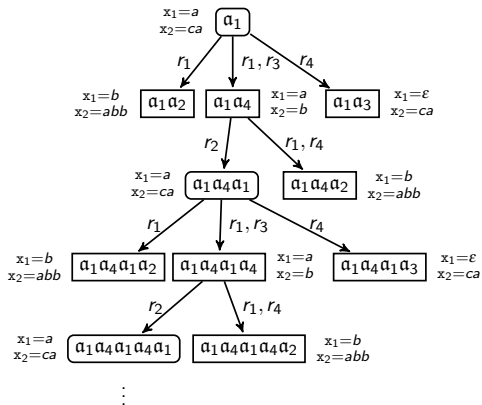
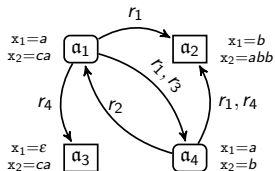
- A concept is in simple form if it is in NNF and all its role paths are of length at most one.
- $\exists rr'. \llbracket S^2y <_{\text{pre}} x \rrbracket$ not in simple form and concepts below in simple form:

$$\exists r. \exists r'. \exists \varepsilon. \llbracket y <_{\text{pre}} x^{\dagger\dagger} \rrbracket \quad \top \sqsubseteq \forall r. \llbracket x = Sx^{\dagger} \rrbracket \quad \top \sqsubseteq \forall r'. \llbracket x^{\dagger} = Sx^{\dagger\dagger} \rrbracket$$

- Given C_0, \mathcal{T} , one can construct in polynomial-time C'_0, \mathcal{T}' in simple form s.t. C_0, \mathcal{T} positive instance of TSAT iff C'_0, \mathcal{T}' positive instance of TSAT.

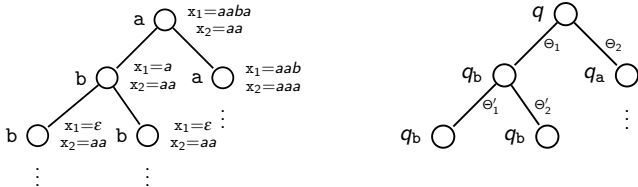


Tree Interpretation Property Needed!



Tree Automata on Strings

(or how to recognize infinite data trees)



- \mathbb{A} accepts infinite trees $\mathbf{t} : [0, d - 1]^* \rightarrow (\Sigma \times (\Sigma^*)^\beta)$.

(ye

- Transitions $(q, a, (\Theta_1, q_1), \dots, (\Theta_d, q_d))$ put constraints on values of current node and children nodes.
- Büchi and Rabin acceptance conditions.
- Can be adapted to many concrete domains and extends similar definitions for the linear case with $d = 1$.

E.g. [Segoufin & Toruńczyk, STACS'11; Kartzow & Weider, arXiv 2015]

From TSAT to Nonemptiness

(or how to apply the standard automata-based approach)

- Technically involved construction following a standard pattern.
 - C_0, \mathcal{T} in simple form positive instance of TSAT iff $L(\mathbb{A}) \neq \emptyset$.
 - Locations are propositionally \mathcal{T} -consistent set of subconcepts.
 - Each role name has dedicated directions in $[1, d]$.
 - Constraints at the level of concepts translated at the level of transitions.
- Postponing the actual problem to the nonemptiness problem.
- Advantages of translating into TCA:
 - Reveals the size of automaton (depending on parameters like number of variables, maximal size for constants, etc.) – important for complexity!
 - Same construction for other concrete domains.

From String Constraints to Integer Constraints

- Intuition: encoding of prefix of strings w and w' by length of common prefix (which is a nonnegative integer).
- $clen(w, w')$: length of longest common prefix btw. w and w' .
E.g. $clen(aba, abbbab) = 2$.
- Properties (I)–(III) are “complete” to *recover string values in a greedy way* ($k = \text{card}(\Sigma)$). [Demri & Deters, JLC 2015]

(I) For $w, w' \in \Sigma^*$, $|w| = clen(w, w) \geq clen(w, w')$.

(II) For all $w_0, w_1, \dots, w_k \in \Sigma^*$ such that

- $clen(w_0, w_1) = \dots = clen(w_0, w_k)$ and,
- for all $i \in [0, k]$, $clen(w_0, w_1) < |w_i|$,

there are $i \neq j \in [1, k]$ such that $clen(w_0, w_1) < clen(w_i, w_j)$.

(III) For all $w_0, w_1, w_2 \in \Sigma^*$,

$clen(w_0, w_1) < clen(w_1, w_2)$ implies $clen(w_0, w_1) = clen(w_0, w_2)$.

Lifting at the Level of Automata

- TCA $\mathbb{A} = (Q, \Sigma, d, \beta, Q_{\text{in}}, \delta, F)$ on \mathcal{D}_Σ translated into $\mathbb{A}' = (Q, \Sigma, d, \beta', Q_{\text{in}}, \delta', F)$ on \mathcal{N} .
- $L(\mathbb{A}) \neq \emptyset$ iff $L(\mathbb{A}') \neq \emptyset$.
- $L(\mathbb{A}') \neq \emptyset$ checked in time

$$R_1(\text{card}(Q) \times \text{card}(\delta') \times \text{MCS}(\mathbb{A}') \times \text{card}(\Sigma) \times R_2(\beta'))^{\mathcal{O}(R_2(\beta') \times R_3(d))}$$

[Demri & Quaas, CONCUR'23]

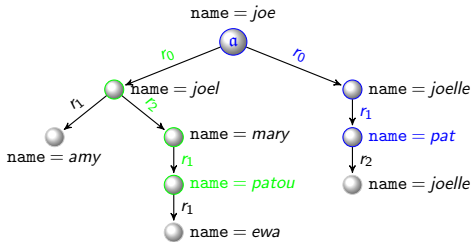
- the R_i 's are polynomials,
- $\text{MCS}(\mathbb{A}')$: maximal size of a constraint in \mathbb{A}' ,
- $\text{MCS}(\mathbb{A}')$ in $(\beta + \text{MCS}(\mathbb{A}) \times \text{card}(\delta) \times d)^{\mathcal{O}(\text{card}(\Sigma)+3)}$,
- β' polynomial in β and in the number of constant strings in \mathbb{A} .

Final Complexity Analysis

C_0, \mathcal{T} is a positive instance iff $L(\mathbb{A}) \neq \emptyset$ and \mathbb{A} satisfies the following quantitative properties:

- Degree d bounded by $\text{size}(C_0, \mathcal{T})$.
- Number of locations in $2^{\mathcal{O}(\text{size}(C_0, \mathcal{T}))}$.
- Number of transitions in $2^{\mathcal{O}(R(\text{size}(C_0, \mathcal{T})))}$ for some polynomial $R(\cdot)$.
- Number of variables β bounded by $\text{size}(C_0, \mathcal{T})$.
- Cardinality of finite alphabet Σ bounded by $2^{\text{size}(C_0, \mathcal{T})}$.
- $\text{MCS}(\mathbb{A})$ quadratic in $\text{size}(C_0, \mathcal{T})$.

Open Problem Related to XPath on Data Trees



$$(\exists r_0 r_1 \cdot \text{name} <_{\text{pre}} \exists r_0 r_2 r_1 \cdot \text{name})$$

- How to handle this extension and characterise its complexity?
EXPTIME-membership ?
- Can we adapt results about XPath on data trees?

See e.g. [Figueira, ToCL 2012]

Concluding Remarks

$$\text{TSAT}(\mathcal{ALCF}^{\mathcal{P}}(\mathcal{D}_{\Sigma})) \longrightarrow \text{TSAT}(\mathcal{ALCF}^{\mathcal{P}}(\mathcal{D}_{\Sigma})) \xrightarrow{\text{in simple form}} \text{NE}(\text{TCA}(\mathcal{D}_{\Sigma})) \longrightarrow \text{NE}(\text{TCA}(\mathbb{N}))$$

- First steps towards taming description logics over strings.
- Automata-based approach with tree constraint automata.
- Reuse or adaptations of several results from literature with new insights to combine them.
- How to extend the results with suffix relation $<_{\text{suf}}$ or regularity constraints?