

# Text Authoring, Knowledge Acquisition and Description Logics

*Marc Dymetman*

Xerox Research Centre Europe

6 chemin de Maupertuis

38240 Meylan

email: marc.dymetman@xrce.xerox.com

## Abstract

We present a principled approach to the problem of connecting a controlled document authoring system with a knowledge base. We start by describing *closed-world authoring* situations, in which the knowledge base is used for constraining the possible documents and orienting the user's selections. Then we move to *open-world authoring* situations in which, additionally, choices made during authoring are echoed back to the knowledge base. In this way the information implicitly encoded in a document becomes explicit in the knowledge base and can be re-exploited for simplifying the authoring of new documents. We show how a Datalog KB is sufficient for the closed-world situation, while a Description Logic KB is better-adapted to the more complex open-world situation. All along, we pay special attention to logically sound solutions and to decidability issues in the different processes.

## Introduction

Recently there has been a surge of interest in *interactive natural language generation systems* (Paris et al., 1995; Power and Scott, 1998; Coch and Chevreau, 2001); such systems rely on a capability of generating a natural language text from an abstract content representation, but — contrary to traditional NLG (Natural Language Generation) systems — this representation is only partially available at the beginning of the text production process; it is then gradually completed by a human author, typically using content-selection menus correlated with regions of the evolving generated text.

One such system, MDA (Multilingual Document Authoring) (citation omitted) is based on a formal specification — using a variant of Definite Clause Grammars (DCGs) (Pereira and Warren, 1980) — of what counts as a valid abstract content representation. The different derivation trees in the grammar correspond to texts with different contents, and at each step of the authoring process the user is asked to make interactive choices on how to expand the current partial derivation tree one step further. There are important analogies between this process and the process of authoring an XML document under the control of a DTD or a Schema, but DCGs are more expressive in terms of the contextual constraints that can be expressed and also are more adapted to the production

of grammatical text.<sup>1</sup>

In published MDA work, all the knowledge about what constitutes a valid document is provided in the grammars, with no clear separation between (1) world knowledge (the fact that a certain pharmaceutical drug contains some molecule makes it dangerous for a certain patient condition) and (2) constraints about document organization (if a certain drug is dangerous for a certain condition, then a warning should be generated at a certain place in the document).

A more principled and modular solution is to leave in the grammar all constraints pertaining to document/textual organization, and to use an external logical theory to express knowledge about the world described by the documents. A document will then be constrained to have a semantic interpretation that is compatible with the external theory.

The aims of this paper are the following.

1. To provide a *formally precise and computationally tractable model* for this approach. The logical theory we will be using will take the form of a Description Logic (DL) knowledge base (Donini et al., 1996); DLs are subsets of FOPC (First-Order Predicate Calculus) which provide a trade-off between expressivity and tractability (in particular decidability) and have recently been given a lot of attention in the knowledge representation community and in activities around the Semantic Web. They are now starting to attract attention in the computational linguistics community as well (Gabsdil et al., 2001; Striegnitz, 2001);
2. To show how this model can be used not only for constraining the document during the authoring process, but also to use the document as a source of new knowledge to be added in a logically sound way to the KB (*knowledge acquisition*);
3. To discuss conditions under which the whole process of authoring is *decidable*.

---

<sup>1</sup>The grammars used in MDA are typically more “semantically” than “syntactically” oriented, and a choice between two alternatives for expanding a nonterminal in the grammar tends to correlate with a clear distinction of meaning in the final text. A given grammar covers a semantically unified class of documents (e.g. employment offers, drug package leaflets, etc.), in a way analogous to the customized XML DTDs used for technical documentation.



content and also to different document instances in the class of documents associated with the grammar. It is then natural to see an abstract tree as a representation of the content of a document belonging to that class.<sup>2</sup>

**Life/death issues** There is one important issue that we did not discuss in the explanation just given, namely how exactly the system determines which choices to propose to the user once he has selected a new nonterminal to be expanded. One possibility is to present him with all the possible names of clauses which are headed by the nonterminal in question (as was done for form(F)), but then it is possible that the author makes a choice that will never lead to a complete valid document.

For instance, let us go back to the point just after the author has chosen tablet as the clause for expanding form(F); at this point the nonterminals on the frontier of the derivation tree are: drug(D), admin(A), comments(D,A), with the constraint df(D,tablet) in the background. Suppose the author next chooses to expand admin(A); if the system was working in a naive fashion, it would then display the choices swallow, chew, and drink. However it is easy to see that drink is in fact ruled out as a choice: any complete document would eventually have to satisfy the constraints df(D,tablet) and da(D,drink), but there is no drug in the database which is compatible with both this form and this administration. We can say that drink is a “dead” choice in this context.

In order to prevent the author from entering a dead-end, what is really needed is for the system to foresee such possible clashes and to present to the author only those choices which may eventually lead to a valid document; in the case at hand, it should present the “live” choices swallow and chew.

**Remark.** When exactly one choice is possible, the system should not even present any choice to the author, but make the only possible expansion decision on its own: authoring should be done automatically at that point. In these cases the authoring mode becomes closer to the classical non-interactive NLG mode, and in the limit, when knowledge-base inferences force all authoring choices, the two modes converge.

**Finitely-parameterized grammars, Datalog, and decidability of life/death** In the current MDA system, the method for determining whether a choice is live or dead is incomplete. This is due to the fact that the nonterminal parameters can be terms of arbitrary complexity (built from variables, constants *and functional symbols*) and then it is easy to simulate with a DCG an arbitrary Prolog program.<sup>3</sup> Determining whether the initial

<sup>2</sup>This abstract tree approach to document content stems from the work of Aarne Ranta on his “Grammatical Framework (GF)” in which he was inspired by the interactive proof editors in a higher-order typed/functional setting such as ALF and COQ in which the user attempts to build a proof of a formula through stepwise top-down refinements of a partial proof (Ranta, 1999). In the present paper the abstract trees can be seen as proofs of an initial goal in a logic programming setting.

<sup>3</sup>Even without the use of auxiliary predicates: a pure Prolog program is equivalent to a DCG generating empty strings.

nonterminal may lead to a complete valid document is then undecidable in general. It is usually possible for the grammar writer to exercise some care in designing the grammars so that life/death problems do not hinder the authoring process in practice, but a principled solution would be preferable.

*In order to tackle this problem, we will be making two fundamental assumptions: (i) the nonterminal parameters in the grammar clauses — as well as the goal arguments in the auxiliary program clauses — are variables or constants; (ii) all variables take their value in the finite set of constants present in the grammar and auxiliary clauses.*

Under these assumptions, we are now dealing with a DCG with *finite-domain parameters* both for its grammar and for its auxiliary predicates components. The auxiliary predicate component is then formally the same as a Datalog database (Ceri et al., 1989), as in our example D1.<sup>4</sup>

We can then see the authoring model as consisting of two components: a finitely parameterized DCG, and a Datalog database.

Now, it is striking that, when working with finite-domain DCGs, not only the auxiliary predicate component, *but also the grammar component*, has formal similarities to a Datalog base: in fact, if one “forgets” in the grammar G1 all the terminal strings, then one obtains a Datalog program DP1:

#### DP1:

```
dfa1:  dfa(D,F,A) ← drug(D), dform(D,F), dadm(D,A).
dform1: dform(D,F) ← form(F), & df(D,F).
dadm1:  dadm(D,A) ← admin(A), comments(D,A).
coms1:  comments(D,A) ← & da(D,A).
coms2:  comments(D,A) ← comments(D,A), comment(D,A).
com1:  comment(D,A).
com2:  comment(diprox,A).
diprox: drug(diprox).
xenor:  drug(xenor).
burpal: drug(burpal).
tablet: form(tablet).
solution: form(solution).
swallow: admin(swallow).
chew:  admin(chew).
drink:  admin(drink).
```

Deciding the productivity of a parameterized nonterminal in the combination G1+D1 is then *formally equivalent to proving it as a program goal* in the combination DP1+D1 (which is itself a global Datalog program), and a derivation in G1 has a one-to-one correspondence to a proof in DP1.

For instance, deciding the productivity of the nonterminal dfa(D,tablet,drink) is equivalent to proving the goal dfa(D,tablet,drink) in the Datalog program DP1+D1: be-

<sup>4</sup>The database D1 only contains facts (Datalog’s EDB), but it could also contain recursively defined predicates (Datalog’s IDB) without impact on the discussion.

cause no such proof can be found, the nonterminal is not productive.

Now, the interest of this translation is that provability of a goal in a Datalog program is not only known to be decidable, but also to be amenable to efficient implementation (Abiteboul et al., 1995).

Consider the situation discussed before, just after the author has chosen the form `tablet`, and at the point where the system needs to present him with a list of choices for `admin(A)`. At that point, the system is confronted with the following question: what are the possible values for `A` such that the following goal:

```
drug(D), admin(A), comments(D,A), df(D,tablet)
```

is satisfiable?

This question can be succinctly represented as the following conjunctive Datalog query:

```
answer(A) ← drug(D), admin(A), com-
ments(D,A), df(D,tablet)
```

for which a number of optimization techniques exist (see (Ceri et al., 1989; Abiteboul et al., 1995)), and which returns as possible values for `A` the set `{swallow, chew}`.<sup>5</sup> The advantage for authoring is clear: at each choice point, the system is capable to return a valid list of choices more efficiently than by applying more naive techniques. It is also worthy of note that some *fundamental issues in authoring are so closely connected with database query optimization*.<sup>6 7</sup>

## Open-world authoring

In an authoring context, some grammatically valid documents will never be authored because they do not correspond to any possible state of affairs. Typically the grammar specifies a much larger set of documents than

<sup>5</sup>In this case, the set of possible values for the parameter `A` coincides with the set of possible values for the names of the expanding clauses for `admin(A)`. In general it is not the case, but it is simple to add a parameter to each nonterminal that indexes its (finitely many) possible expanding clauses.

<sup>6</sup>A DCG is nothing else than a context-free grammar with parameterized nonterminals and a unification mechanism between the parameters. Because of the analogy between DTD/Schemas and CFGs, it seems likely that the same approach could be useful for extending XML-based authoring through the use of finite-domain parameters and unification.

<sup>7</sup>The fact that the program `DP1` is equivalent to `G1` as far as *non-terminal productivity is concerned* does not mean that the two objects are equivalent for authoring purposes. The grammar associates different texts with different derivations of the *same* ground nonterminal (for instance, there are an infinite number of texts produced by `COMMENTS(diprox,tablet)`, corresponding to different combinations of `coms1`, `coms2`, `com1`, `com2`), whereas the program is of interest to us here not in the different proofs of a given ground goal, but in the fact that this goal is provable or not. Note that the clause of name `COMS2` can be eliminated from the program `DP1` without changing its interpretation (because in order to prove `COMMENTS(D,A)` it requires a proof of `COMMENTS(D,A)`), but making the program non-recursive and therefore simplifying the check for productivity; eliminating the same clause from `G1` would however completely change the meaning of the grammar.

the ones which are actually possible. If this were not the case, then an author would not have to take the trouble to direct the production process by making content choices that he alone can make. That is to say, a document which has actually been authored conveys more meaning than just stating “I am a valid document relative to the specification”. However, in a closed-world environment as we have been discussing until now, that additional meaning has no explicit counterpart in the knowledge-base; it is only represented implicitly in the abstract content tree, in a form which is not perspicuous and would be difficult to re-use for the authoring of other documents or to share with other processes.<sup>8</sup>

In a closed-world context, the KB constraints which are tested during the authoring process are completely passive: they are seen purely as validity checks against the knowledge base.

By contrast, open-world authoring sees the KB constraints not only as checks, but also as conditions on the world being described. When authoring a document, the author is not neutrally picking out one of the documents valid relative to the KB, but asserting that the constraints *do* hold of the actual world.<sup>9</sup>

Let us illustrate this idea. We are now viewing the formal specification of valid documents as consisting, as before, of a grammar of the type previously described (we will take again the grammar `G1`), but instead of a Datalog database, we are now using an *informationally incomplete* description logic knowledge base `KB1`:<sup>10</sup>

**KB1:**

TBOX:

```
TabletDrugs = ∃df.{tablet}
SolutionDrugs = ∃df.{solution}
SwallowDrugs = ∃da.{swallow}
ChewDrugs = ∃da.{chew}
DrinkDrugs = ∃da.{drink}
```

```
Drugs = TabletDrugs ⊔ SolutionDrugs
Drugs = SwallowDrugs ⊔ ChewDrugs ⊔ DrinkDrugs
```

```
TabletDrugs = SwallowDrugs ⊔ ChewDrugs
SolutionDrugs = DrinkDrugs
```

ABOX :

```
df(burpal,solution)
da(burpal,drink)
```

This knowledge-base is written using a certain number of DL constructors — existential quantification, concept

<sup>8</sup>Note an analogy here with the Semantic Web perspective: tags used in XML documents may convey implicit semantic information, but in order to make this information sharable, it had better be represented explicitly in some formal knowledge representation.

<sup>9</sup>In the language of pragmatics, the author is then performing a *speech act* by committing to the “truth” of the document.

<sup>10</sup>An introduction to DLs would take us too far afield; let’s just say that there is a whole family of DLs, which differ by the logical constructors they allow, and that most can be seen as decidable fragments of first-order logic. An accessible recent introduction to DLs is available at <http://www.cs.man.ac.uk/horrocks/Slides/leipzig-jun-01.pdf>.

enumeration, disjoint union (an abbreviation:  $A = B \uplus C$  can be replaced by the two constraints  $A = B \sqcup C$  and  $B \sqcap C = \perp$ , and  $B \uplus C \uplus D$  is an abbreviation for  $(B \uplus C) \uplus D$ ) —, and we are assuming the unique name convention (all named individuals are different). The constructors which are used place the knowledge base in the class *ALCO* (Donini et al., 1996).

The TBOX can be glossed in the following way. The TabletDrugs are those drugs  $D$  for which  $df(D, \text{tablet})$ , the SolutionDrugs those drugs for which  $df(D, \text{solution})$ , ..., the DrinkDrugs those drugs for which  $da(D, \text{drink})$ . The drugs can come in either one of the two forms: tablet and solution, and in either one of the three administrations swallow, chew and drink. Finally TabletDrugs are either swallow drugs or chew drugs, whereas SolutionDrugs are always drink drugs. The ABOX says what we already know about the form and administration of Burpal.

The list of relations in  $D1$  is compatible with  $KB1$ : indeed it is easy to see that one can obtain a model of  $KB1$  by taking the relations of  $D1$  along with the facts:

```
diprox: TabletDrugs
xenor : TabletDrugs
burpal : SolutionDrugs
diprox: SwallowDrugs
xenor : ChewDrugs
burpal : DrinkDrugs
```

In a certain sense the TBOX of  $KB1$  can be seen as a conceptual schema for the database  $D1$ , which states certain general relations about the forms and administrations of drugs, or about the uniqueness of form and administration for a drug, but which does not say how many drugs there are or what are the properties of these drugs.

**Valid abstract trees and incomplete KBs** Let us return to our authoring example in this new context. We now associate grammar  $G1$  with  $KB1$  instead of  $DB1$ . We then make the assumption that all constant parameters appearing in the grammar (diprox, xenor, burpal, tablet, etc.) are to be considered distinct named individuals for the KB, and that the constraint relations (da, df) are all unary or binary and correspond to concepts or roles in the KB.

Let's now look again at the abstract tree  $AT1$ :

```
dfa1(diprox, dform1(tablet), dadm1(swallow, coms1))
```

This abstract tree is valid *relative to  $G1$*  (it corresponds to a possible complete derivation) but it is not necessarily valid relative to the combination  $\langle G1, KB1 \rangle$ ; this notion is defined in the following way: because the abstract tree uniquely determines the set of rules which have been used for building the derivation, it also uniquely determines a set of associated KB constraints; thus  $AT1$  is associated with the set of constraints:  $\{df(diprox, \text{tablet}), da(diprox, \text{swallow})\}$ .

Now we say that  $AT1$  is valid relative to the combination  $\langle G1, KB1 \rangle$  if and only if it is both valid relative to  $G1$  and if its associated set of constraints is

compatible with  $KB1$ . In other words we need to show that *the addition of the two constraints  $df(diprox, \text{tablet}), da(diprox, \text{swallow})$  to the ABOX still leads to a satisfiable knowledge base*. This can be shown by exhibiting a model as we did a few paragraphs ago, and therefore  $AT1$  is a valid abstract tree relative to  $\langle G1, KB1 \rangle$ .

The informal reasoning by which we just showed the satisfiability of  $KB1$  extended with the two relations can also be established by a computational proof, due to the decidability of KB-consistency checking in *ALCO* (Donini et al., 1996).

**Open- vs. closed-world authoring, satisfiability vs. deducibility** Note that validity of an abstract tree in the *open-world authoring* context involves the *satisfiability* of a conjunction of constraints relative to the knowledge base, whereas the notion of validity of an abstract tree in the *closed-world authoring* context involves the dual notion of *deducibility* of a conjunction of constraints relative to the knowledge-base (in the Datalog context, being true in the minimal Herbrand model is the same as being deducible from the Horn clauses constituting the base).

**Decidability of the authoring process** In order to illustrate the process, let's go back to the point in the authoring after all obligatory expansions of  $dfa(D, F, A)$  have been made, where the frontier of the derivation tree is  $drug(D)$ ,  $form(F)$ ,  $admin(A)$ ,  $comments(D, A)$ , and where the user has chosen to expand  $form(F)$ . There are apparently two possible expansions: the clauses with names *tablet* and *solution*. Before presenting these choices to the user, the system must check that they are live, namely, as before, that they may lead to a complete valid document.

Choosing the *tablet* expansion leads to the derivation frontier  $drug(D)$ ,  $admin(A)$ ,  $comments(D, A)$  with constraint  $df(D, \text{tablet})$ . In order to decide whether the frontier is live, the system needs to enumerate possible complete derivations of this frontier until it finds one that is satisfiable relative to  $KB1$  and then return a positive answer, and if it does not find one, it should return a negative answer. In principle, the enumeration could never stop, but because of the finite parameter condition on the grammar, the system has only to enumerate a finite number of trees; this is because if a derivation tree is of the form  $S(\dots A1(\dots A2(\dots) \dots) \dots)$  where  $S$  is a ground instantiation of the initial nonterminal and  $A1$  and  $A2$  are the *same* ground instantiation of a nonterminal ("repetitive derivation"), then the satisfiability of  $S(\dots A1(\dots A2(\dots) \dots) \dots)$  relative to  $KB1$  implies the satisfiability of  $S(\dots A2(\dots) \dots)$ : a model of the larger derivation tree is again a model of the smaller derivation tree. *This means that when checking life/death we do not ever need to consider a repetitive derivation during the enumeration of derivations*. In particular, because we are dealing with a finite parameter domain, the derivations that we need to consider have a bounded depth (otherwise we would necessarily encounter repetitive situations), and the decidability of the process follows.<sup>11</sup>

<sup>11</sup>The same reasoning could be made for proving decidability in the

In the case of choosing tablet, the abstract tree AT1 is enumerated at some point in the process, and its satisfiability relative to KB1 can be decidably checked: tablet is then shown to be a live authoring choice. The same process shows solution to be live.

Now, let's go to the point where, after having chosen tablet, the author decides to select an expansion for admin(A). The derivation frontier is then drug(D), admin(A), comments(D,A), with the constraint df(D,tablet), and the apparently possible expansions are swallow, chew, and drink. Both swallow and chew can be seen to be live by a similar reasoning as before. In the case of drink, we have to check whether the sequence drug(D), comments(D,drink), with the constraint df(D,tablet) is live. Let's choose to expand comments(D,drink) first. The expansion coms2 leads to a repetitive situation (comments(D,drink) is above comments(D,drink) in the derivation path.) and is therefore discarded; the expansion coms1 leads to the frontier drug(D), with the constraints df(D,tablet) and da(D,drink). However the two constraints cannot be simultaneously satisfied in KB1; This can be shown computationally by using the satisfiability check in KB1, but also by the following informal reasoning: df(D,tablet) and da(D,drink) imply that D is both in TabletDrugs and in DrinkDrugs; by the second fact it is in SolutionDrugs, but SolutionDrugs and TabletDrugs have an empty intersection. Thus all expansions of comments lead to invalidity; hence drink is not a live choice.

**Open-World authoring and hybrid knowledge bases** The process that we have just described for finding live selections, although decidable, is clearly not optimized. In the case of closed-world authoring that we discussed at the beginning of this paper, we said that, from the point of view of detecting life/death situations, a Datalog program such as DP1 could be used in place of the grammar G1, and that the combination of DP1 + D1 could be treated as a global Datalog database to which standard query optimization techniques could be applied. Is there some comparable possibility here? A clue comes from the area of *hybrid knowledge bases* in the description logic community. Some researchers have shown that by associating Description Logics with Datalog one can significantly increase the expressive power of both formalisms, which have a nice complementarity (recursive definitions can be easily expressed in Datalog, but not in DLs; partial knowledge can be easily expressed in DLs, but not in Datalog) (Levy and Rousset, 1996; Donini et al., 1998). The open-authoring approach we propose has strong connections with these hybrid knowledge-bases (citation omitted) and it seems likely that optimization techniques from that area may be transferred to our problem.

**Light semantics and knowledge acquisition** Let's step back and reconsider the rationale behind open-world authoring. We are considering a situation in which there is an "actual world" which is not completely known ei-

closed-world case, instead of appealing there to the decidability of Datalog queries.

ther to the knowledge base or to the author; however both the KB and the author are supposed to have correct partial knowledge about that world.

The system presents the author with a collection of documents which, *from its point of view*, are compatible with what *it* knows about the actual world. Among these documents, the author picks (during the authoring process) one document that, *from his point of view*, is compatible with what *he* knows about the actual world.

So the author is not passively exploring the space of document considered possible by the system (although that could certainly be a nonstandard mode of operation if the author takes a developer's hat and wants to see what the system believes is possible), but is actively committing to certain facts about the world.

What are these facts? What the author is producing is an abstract content tree, which corresponds to a *completely specific choice of expansion rules* for the nonterminals of the grammar. This means that the abstract tree completely determines a set of associated ground KB relations. For instance AT1 determines the set {df(diprox,tablet), da(diprox,swallow)}. These are the facts that the author asserts to be true in the actual world.

**Light semantics.** Such facts are aspects of the document content that the document "exports" to the knowledge base and thereby makes formally explicit. They provide what we shall call a *light semantics* for the document. In terms of light semantics, if we were to build a standard logical form for the whole document, for instance for AT1, that logical form would simply be the conjunction of the associated asserted facts  $df(diprox,tablet) \wedge da(diprox,swallow)$ . Light semantics does not attempt to model the whole semantics of the document (for instance, in our example, there is no explicit logical counterpart to the different choices for the comment nonterminal), but focuses instead on modeling those parts of the document semantics that can be tractably handled both by the knowledge representation component and by the authoring process.<sup>12</sup>

**Knowledge acquisition.** Once the author has committed to a document, he has revealed a certain number of facts that he knows about the actual world and that the

---

<sup>12</sup>When working in a more powerful framework for logical forms, such as Montague semantics, the interpretation of a document may depend in non-monotonic ways on the interpretations of its parts, as in negated contexts: "it is not the case that ..." or in opaque contexts: "John believes that ...". Predicting at authoring time which selections are live relative to such a knowledge representation framework, while possible in principle, seems to be a difficult research question. Another (orthogonal) argument in favor of light semantics is the fact that if we consider the communicative role of a document inside a predefined class of documents, then there is no point in formally representing those parts of a document that are not contrastive between two documents in the class; for instance, there is no need to analyze the sentence "Always ask your doctor's advice in case of doubt" in any semantic detail if it appears in all documents of the class: these semantic details are irrelevant to the informational content of the document as opposed to other documents of the class. A thorough discussion of this point, connected to considerations of information theory, would bring us outside the scope of this paper.

KB possibly did not “know”. These facts (in our example:  $df(diprox,tablet)$  and  $da(diprox,swallow)$ ) can then be added to the ABOX of the knowledge base, and can be used either for their own sake (knowledge acquisition) or in order to constrain the authoring of a new document.

So after the authoring of AT1, the ABOX of KB1 becomes:

ABOX :  
 $df(burpal,solution)$   
 $da(burpal,drink)$   
 $df(diprox,tablet)$   
 $da(diprox,swallow)$

Suppose now the user authors a new document, first making a selection for drug(D), and choosing  $diprox$ . Then the KB “knows” that  $tablet$  is the only choice for F and  $swallow$  the only choice for A. Indeed they are possible choices (because  $df(diprox,tablet)$  and  $da(diprox,swallow)$  are in the ABOX of the KB), but are also the *only* choices, for  $diprox$  is now known to be in  $TabletDrugs$  and in  $SwallowDrugs$ ; it can therefore not be in  $SolutionDrugs$  or in  $ChewDrugs$  or in  $DrinkDrugs$ , which means that none of the facts  $df(diprox,solution)$ ,  $da(diprox,chew)$  or  $da(diprox,drink)$  may hold. After the author’s choice of  $diprox$ , the derivation frontier is  $form(F)$ ,  $admin(A)$ ,  $comments(diprox,A)$  with the constraint  $df(diprox,F)$ . The author then chooses to expand  $form(F)$ , and the system notices that the only live choice is  $tablet$ , and performs this expansion without asking the user. The frontier is now  $admin(A)$ ,  $comments(diprox,A)$ , with the constraint  $df(diprox,tablet)$ . Now the user can choose to expand  $admin(A)$ , and the only live choice is  $swallow$ . At that point the frontier is  $comments(diprox,swallow)$  with the constraint  $df(diprox,tablet)$ . The author can then make choices for  $comments(diprox,swallow)$  that lead to zero or several instances of  $comment(diprox,swallow)$ . At a certain point he will choose the nonrecursive expansion  $com1$ , which will lead to an empty frontier, with the constraints  $df(diprox,tablet)$  and  $da(diprox,swallow)$ .

We could obviously suppose here that rather than waiting for the user to point to the nonterminal he wants to expand next before finding the live choices for this nonterminal, the system could find all the live choices for all nonterminals on the frontier beforehand, and do the obligatory expansions without any input from the user, but at a slightly higher computational cost. In this way, after the initial choice of  $diprox$  as the drug, the other steps of the authoring process would be done automatically, apart from the choice of how many (and which) comments to make, which would still remain the responsibility of the author.

## Conclusion

In the course of the paper we have defined different notions such as *live-death* issues in authoring processes, *closed-world* versus *open-world* authoring, and *light document semantics*. We have presented a formal ap-

proach to closed-world authoring that shows a correspondence between life-death problems and conjunctive Datalog queries, as well as a formal approach to open-world document authoring based on Description Logics. We have also sketched proofs of decidability for life-death issues in these different processes. Finally we have shown how an open-world authoring context can be used for supporting a novel form of *knowledge acquisition*.

## Acknowledgments

Many thanks to Jean-Marc Andréoli, Caroline Brun, Éric Fanchon, Pierre Isabelle, Aaron Kaplan, Aurélien Max, and Sylvain Pogodalla for discussions and comments, and to the anonymous reviewers for suggestions on improving the paper.

## References

- Serge Abiteboul, Richard Hull, and Victor Vianu. 1995. *Foundations of Databases*. Addison-Wesley.
- S. Ceri, G. Gottlob, and L. Tanca. 1989. *Logic Programming and Databases*. Springer-Verlag.
- José Coch and Karine Chevreau. 2001. Interactive multilingual generation. In A. Gelbukh, editor, *Computational Linguistics and Intelligent Text Processing, LNCS 2004*. Springer.
- Francesco M. Donini, Maurizio Lenzerini, Daniele Nardi, and Andrea Schaerf. 1996. Reasoning in description logics. In Gerhard Brewka, editor, *Principles of Knowledge Representation*, pages 191–236. CSLI Publications, Stanford, California.
- Francesco M. Donini, Maurizio Lenzerini, Daniele Nardi, and Andrea Schaerf. 1998. AL-log: Integrating datalog and description logics. *Journal of Intelligent Information Systems*, 10(3):227–252.
- Marc Dymetman. 2002. Document content authoring and hybrid knowledge bases. In *Proceedings of the 9th International Workshop on Knowledge Representation meets Databases (KRDB-2002)*, Toulouse, April.
- Malte Gabsdil, Alexander Koller, and Kristina Striegnitz. 2001. Building a text adventure on description logic. In *Proceedings of KI-2001 Workshop on Applications of Description Logics*, Vienna.
- Alon Y. Levy and Marie-Christine Rousset. 1996. CARIN: A representation language combining horn rules and description logics. In *European Conference on Artificial Intelligence*, pages 323–327.
- Cécile Paris, Keith Vander Linden, Markus Fischer, Anthony Hartley, Lyn Pemberton, Richard Power, and Donia Scott. 1995. A Support Tool for Writing Multilingual Instructions. In *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI) 1995*, pages 1398–1404, Montréal, Canada.
- F. Pereira and D. Warren. 1980. Definite clauses for language analysis. *Artificial Intelligence*, 13:231 – 278, 1980.
- Richard Power and Donia Scott. 1998. Multilingual authoring using feedback texts. In *COLING-ACL*, pages 1053–1059.
- Aarne Ranta. 1999—. Grammatical framework work page. [www.cs.chalmers.se/~aarne/GF/pub/work-index/index.html](http://www.cs.chalmers.se/~aarne/GF/pub/work-index/index.html).
- K. Striegnitz. 2001. Model checking for contextual reasoning in nlg. ICOS-3. Inference in Computational Semantics Workshop. Siena.