

FOUR-VALUED KNOWLEDGE AUGMENTATION FOR STRUCTURED DOCUMENT RETRIEVAL

M. LALMAS and T. ROLLEKE

*Department of Computer Science, Queen Mary University of London, United Kingdom
{mounia,thor}@dcs.qmul.ac.uk*

Received (received date)

Revised (revised date)

Structured documents are composed of objects with a content and a logical structure. The effective retrieval of structured documents requires models that provide for a content-based retrieval of objects that takes into account their logical structure, so that the relevance of an object is not solely based on its content, but also on the logical structure among objects. This paper proposes a formal model for representing structured documents where the content of an object is viewed as the knowledge contained in that object, and the logical structure among objects is captured by a process of knowledge augmentation: the knowledge contained in an object is augmented with that of its structurally related objects. The knowledge augmentation process takes into account the fact that knowledge can be incomplete and become inconsistent.

Keywords: structured document, knowledge augmentation, formal model

1. Introduction

With the widespread development of structured document repositories (e.g. CD-ROM, the Internet and digital libraries) and the rapid adoption of the XML markup language, there is more scope and need to exploit the structural knowledge of documents for the purpose of their retrieval.^{1,2,3,4,5,6,7}

In this paper, we are concerned with structured documents composed of *objects* with a content and a logical structure.^{8,9,10} These objects correspond to the document components. The content refers to the content of objects. For example, the content of an object can be described by a term “sailing”. The logical structure refers to the way structured documents are organised. For example, a structured document may be composed of objects such as sections and figures, which can themselves be composed of other objects (e.g. a section and its subsections). The *root* object, which is unique, embodies the whole document. *Atomic* objects are document components that are not composed of other components. All other objects are referred to as *inner* objects.

A number of approaches have been developed to specifically deal with structured document retrieval. They can be classified into three groups. Passage retrieval approaches aim at retrieving documents based on the most relevant passage(s) in documents.^{11,12,13} Data modelling approaches aim at developing data models for representing and querying with respect to the content and structure of documents.^{14,15,16} Aggregation-based approaches represent or estimate the relevance of objects based on the aggregation of the representation or estimated relevance of their own content and the representation or estimated relevance of their structurally related objects.^{10,3,17,18,19,20,5} All three types of approaches have highlighted that considering structure and content knowledge can improve the retrieval effectiveness of structured documents.

In this paper, we follow the aggregated-based approach to structured document retrieval, which allows the retrieval of objects of varying granularity as a result to a query.^{3,21,18,20,22} That is, retrieved objects consist of entry points to a same document: an entry to the root object when the entire document is estimated relevant to a query, an entry to an atomic object when only that object is estimated relevant to the query, or an entry to an inner object when that object and its structurally related objects are estimated relevant to the query.

This paper proposes a formal model for representing structured documents that allows the retrieval of objects of varying granularity.¹ The model exploits the content and the logical structure among objects to arrive at a representation of the structured document. The model developed in this paper is based on the work of Chiaramella et al,³ where the basis of an aggregated-based approach to structured document retrieval was proposed.

The content of an object is viewed as the *knowledge* contained in that object, and the structure among objects is captured by a process of *knowledge augmentation*: the knowledge contained in an object is *augmented* with the knowledge contained in its structurally related objects. The knowledge augmentation process is a combination of knowledge specifically defined to provide for a representation of the object that is based on its own content and that of its structurally related objects. For instance, a non-atomic object composed of an object about “sailing” and a second object about “boat” can be considered as being about “sailing and boat”. The knowledge contained in the non-atomic objects (its content description) is augmented with that of its related objects. It is this representation that can then be used to estimate the relevance of the object to a query. The estimate takes into account that non-atomic objects are structurally related to other objects. In our example, the non-atomic object should be ranked higher than its related objects for a query about “sailing and boats”.

In this paper, we are concerned with the representation of structured documents, in particular the knowledge augmentation process. Our proposed model formally

¹This paper is an extended version of the paper entitled “Four-valued knowledge augmentation for representing structured documents”, appearing in the 13th International Symposium on Methodologies for Intelligent Systems (ISMIS02), pp 158-166, Lyon, France, June 2002.

expresses this process.

Our model formally takes into account the fact that knowledge can be incomplete and become inconsistent. Incompleteness is due to the fact that knowledge that is not explicit in an object should not be considered false. For instance, an object composed of one object about “sailing” and a second object where nothing is said about “sailing” (incomplete description) should be considered as being about “sailing”. Inconsistent arises when two objects upon which the augmentation process is based contain knowledge that is contradictory. For example, an object composed of one object about “sailing” and a second object about “not sailing” should not be considered as being about “sailing” nor “not sailing” since there is contradictory evidence.

The model is based on the definition of modal operators, referred to as *knowledge modal operators*, one for each of the objects forming a structured document. The knowledge modal operator associated with an object is used to model the knowledge contained in the object, which can be of two types: content knowledge (i.e. propositions describing the information content of the object) and structural knowledge (i.e. the fact that the object is structurally related to other objects). The semantics of the knowledge modal operators is based on an interpretation structure and an interpretation function defined upon four truth values, possible worlds, accessibility relations and truth value assignment functions.^{23,24} Possible worlds and truth value assignment functions reflect the knowledge contained in objects. The accessibility relations captures structural relationships between objects. Four truth values are used to capture incompleteness and inconsistency.²⁵

The knowledge augmentation process is formalised through the definition of a different set of modal operators, referred to as *augmented knowledge modal operators*, one for each of the objects forming a structured document. An augmented knowledge operator is used to model the content knowledge of an object augmented with that of its structurally related objects. The semantics of the augmented knowledge modal operators is defined upon the interpretation structure and interpretation function defining the semantics of the knowledge modal operators.

We also show that the knowledge augmentation process can be described using the framework of Fagin et al,²⁴ for modelling combination of knowledge. This is achieved through the definition of *G-world trees* that formalise a graphical representation of the set of possible worlds and accessibility relations associated with the representation of objects. Special truth value assignment functions, referred to as *augmented truth value assignment functions*, defined upon G-world trees are used to characterise the augmented knowledge. We show that the two approaches are equivalent, that is, the same representation of objects is obtained after the knowledge augmentation process.

The outline of this paper is as follows. Section 2 describes the model for representing content and structural knowledge of objects. Sections 3 and 4 present the two approaches for describing the knowledge augmentation process. Section 5 shows that the two approaches are equivalent. Section 6 concludes and discusses

future work.

2. Representing content and structural knowledge

The syntax² used to characterise the content and structural knowledge of a structured document is given in Figure 1. A structured document is described by a *program*, which is a set of *clauses* where a clause is either a *fact* or a *context*. A *context* consists of an object identified by an object name, in which a program is nested. The object name in a context clause is referred to as a *context name*. We use the context name “this” to refer to the context name associated with the global program (i.e. representing the database of structured documents). The set of object names is called \mathcal{C} .

A *fact* is a *proposition*, or a proposition preceded by “not” or *truth-list*. We call Φ the set of propositions. *truth-list* represents the four truth values: “1” means *true*, “0/1” means *false*, “0/0/1” means *inconsistent*, and “0/0/0/1” means *unknown*. In this paper, “1 proposition” and “0/1 proposition” are the same as “proposition” and “not proposition”, respectively. We use both “0/1” and “not” because the proposed representation is currently being extended to allow for a probabilistic-based representation of objects, thus capturing the uncertainty inherent to the information retrieval process.

program	::=	clause clause program
clause	::=	fact context
fact	::=	proposition 'not' proposition truth-list proposition
proposition	::=	NAME
truth-list	::=	'1' '0/1' '0/0/1' '0/0/0/1'
context	::=	NAME '[' program ']'
NAME	::=	[a-z][a-zA-Z0-9]*

Fig. 1. Syntax of programs

Let *doc* (e.g. a document) and *sec* (e.g. a section) be context names in \mathcal{C} and “sailing” be a proposition in Φ . A context clause is called an *atomic context* clause when it contains only facts, but no context clauses. For instance, “doc[sailing]” is an atomic context clause and expresses that the content of the object modelled by the context name *doc* is described by the proposition “sailing”. Otherwise, a context clause is called a *structured context* clause. For instance, “doc[sec[sailing]]” is a structured context clause and expresses that the object represented by *doc* is composed of the object represented by *sec*, and the content of the object modelled by *sec* is described by the proposition “sailing”. We refer to *sec* as a *subcontext* of

²This is not the full syntax. Our full formalism provides for the definitions of sentences built upon the conjunction and disjunction of propositions/sentences, rules, queries, etc. Only the part relevant to the knowledge augmentation process is described in this paper.

doc, and doc as *the supercontext* of sec. A supercontext can have several subcontexts, whereas a subcontext has exactly one supercontext. Atomic context clauses characterise atomic objects, whereas structured context clauses characterise non-atomic objects (i.e. inner and root objects). Finally, a function $sub : \mathcal{C} \rightarrow \wp(\mathcal{C})$ yields the set of subcontexts of a supercontext for a given program.

Figure 2 is a commented example of a program (the comments appear after %). The structured document is composed of two objects represented by the structured context doc and the atomic context sec. doc is a subcontext of the global context, denoted “this”, representing the whole database (e.g. the set of all structured documents). sec is a subcontext of doc. The context sec has two propositions “not boats” and “sailing”. Finally, $sub(\text{this}) = \{\text{doc}\}$ and $sub(\text{doc}) = \{\text{sec}\}$.

```

doc[   % structured context
  sec[ % atomic context
    not boats      % same as 0/1 boats
    sailing        % same as 1 sailing
  ] % end of context sec
] % end of context doc

```

Fig. 2. Example of a program

The semantics of programs is defined upon an *interpretation structure* denoted M and an *interpretation function* denoted \models that assigns truth values to programs, clauses, contexts, facts and propositions. We consider contexts as “agents” that possess *knowledge* concerning their content and their structure and use a Kripke structure to define the semantics of programs. Different from their work, we use the terminology “context” instead of “agent” and we consider four truth values. We first define an interpretation structure with respect to a set of propositions Φ and a set of context names \mathcal{C} .

Definition 1 (Interpretation structure M): An interpretation structure M for a set Φ of propositions and a set \mathcal{C} of context names is a tuple $M = (W, \pi, \mathcal{R})$ where

- $W := \{w_d : d \in \mathcal{C}\} \cup \{w_{slot}\}$ is a finite set of possible worlds. For each context name, a possible world is defined. The set includes the possible worlds w_{this} for the global context “this” and w_{slot} to model the access to the “this” context.
- $\pi : W \rightarrow (\Phi \rightarrow \{\text{true}, \text{false}, \text{inconsistent}, \text{unknown}\})$ is a function on W that yields a truth value assignment for all world w in W , which is a function

$\pi(w) : \Phi \rightarrow \{\text{true}, \text{false}, \text{inconsistent}, \text{unknown}\}$ that assigns a truth value to each proposition in Φ , where the four truth values are defined as $\text{true} := \{t\}$, $\text{false} := \{f\}$, $\text{inconsistent} := \{t, f\}$, and $\text{unknown} := \{\}$.

- $\mathcal{R} := \{R_d : d \in \mathcal{C}\}$ is a finite set of binary relations on W , called accessibility relations, one for each context name. For any supercontext d and subcontext s , i.e. $s \in \text{sub}(d)$, $\{(w_d, w_s)\} \in R_s$ where R_s is the accessibility relation associated to s , and w_d and w_s are the possible worlds associated with d and s , respectively. We say that context s accesses or reaches world w_s from world w_d . The accessibility relation for the global context “this” is defined as $R_{\text{this}} := \{(w_{\text{slot}}, w_{\text{this}})\}$.
- The function $R_d(w) := \{w' \mid (w, w') \in R_d\}$ for $w \in W$ and $R_d \in \mathcal{R}$ yields the set of worlds that can be reached from world w by the context d .

Consider the program “doc[sec[sailing]]”. doc and sec are context names and their respective possible worlds are w_{doc} and w_{sec} . To represent that the proposition “sailing” appears within sec but neither in doc or the global context “this”, the truth values assigned to sailing are as follows: $\pi(w_{\text{sec}})(\text{sailing}) = \text{true}$, $\pi(w_{\text{doc}})(\text{sailing}) = \pi(w_{\text{this}})(\text{sailing}) = \pi(w_{\text{slot}})(\text{sailing}) = \text{unknown}$. *unknown* captures incompleteness since it is not explicitly stated that sailing does not describe the content of doc and the global context. The logical structure is represented through the accessibility relations $R_{\text{sec}} = \{(w_{\text{doc}}, w_{\text{sec}})\}$, $R_{\text{doc}} = \{(w_{\text{this}}, w_{\text{doc}})\}$ and $R_{\text{this}} = \{(w_{\text{slot}}, w_{\text{this}})\}$.

The semantics of the content and structural knowledge of objects is based upon considering context names as modal operators, referred to as *knowledge modal operators*. The atomic context clause “doc[sailing]” becomes interpreted as “doc knows sailing” and captures the content knowledge of the corresponding object. The structured context clause “doc[sec[sailing]]” becomes interpreted as “doc knows that sec knows sailing” and captures the structural knowledge of the object represented by doc. With this interpretation in mind, we define the interpretation function \models that assigns truth values to propositions, facts, atomic and structured context clauses, and programs with respect to the interpretation structure defined above.

Definition 2 (Interpretation function \models): Let $M = (W, \pi, \mathcal{R})$ be a interpretation structure as defined in Definition 1. Let φ be a proposition in Φ and w a world in W . Let d be a context name in \mathcal{C} and $R_d \in \mathcal{R}$ its corresponding accessibility relation. Let s be a context name in \mathcal{C} . The interpretation of facts is defined as follows:

$$\begin{aligned}
 (M, w) \models 1 \varphi & : \iff \pi(w)(\varphi) = \text{true} \\
 (M, w) \models 0/1 \varphi & : \iff \pi(w)(\varphi) = \text{false} \\
 (M, w) \models 0/0/1 \varphi & : \iff \pi(w)(\varphi) = \text{inconsistent} \\
 (M, w) \models 0/0/0/1 \varphi & : \iff \pi(w)(\varphi) = \text{unknown}
 \end{aligned}$$

The interpretation of atomic context clauses is defined as follows:

$$\begin{aligned}
(M, w) \models d[1 \varphi] & : \iff \forall w' \in R_d(w) : (M, w') \models 1 \varphi \\
(M, w) \models d[0/1 \varphi] & : \iff \forall w' \in R_d(w) : (M, w') \models 0/1 \varphi \\
(M, w) \models d[0/0/1 \varphi] & : \iff \forall w' \in R_d(w) : (M, w') \models 0/0/1 \varphi \\
(M, w) \models d[0/0/0/1 \varphi] & : \iff \forall w' \in R_d(w) : (M, w') \models 0/0/0/1 \varphi
\end{aligned}$$

The interpretation of structured context clauses is defined as follows:

$$(M, w) \models d[s[\text{fact}]] : \iff \forall w' \in R_d(w) : (M, w') \models s[\text{fact}]$$

The interpretation of a program within a context d is defined as follows, where program^* is the set of clauses given in the program:

$$(M, w) \models d[\text{program}] \iff \forall \text{clause} \in \text{program}^* : (M, w) \models d[\text{clause}]$$

An interpretation structure M is called a model of a program P iff $\text{this}[P]$ is true in all worlds with respect to M . We use “valid” for “true in all worlds”. The notation $M \models \text{this}[P]$ reads “ $\text{this}[P]$ is valid”.

Returning to our example, for an interpretation M to be a model of the program “ $\text{this}[\text{doc}[\text{sec}[\text{sailing}]]]$ ”, we need to show that for all worlds $w \in W$,

$$(M, w) \models \text{this}[\text{doc}[\text{sec}[\text{sailing}]]]$$

These worlds are $w_{slot}, w_{this}, w_{doc}$ and w_{sec} . We have $\pi(w)(\text{sailing}) = \text{true}$ in w_{sec} and unknown in all other worlds. We start with w_{slot} . For $(M, w_{slot}) \models \text{this}[\text{doc}[\text{sec}[\text{sailing}]]]$ to hold, we must verify that $(M, w_{this}) \models \text{doc}[\text{sec}[\text{sailing}]]$ holds since w_{this} is the only world accessible from w_{slot} with respect to R_{this} . For $(M, w_{this}) \models \text{doc}[\text{sec}[\text{sailing}]]$ to hold, we must verify that $(M, w_{doc}) \models \text{sec}[\text{sailing}]$ holds, which holds if $(M, w_{sec}) \models \text{sailing}$, which is holding from the truth value assignment function $\pi(w_{sec})$. In all other worlds w_{this}, w_{doc} and w_{sec} , the context “this” cannot reach another world (i.e., through the accessibility relation R_{this}), hence “ $\text{this}[\text{doc}[\text{sec}[\text{sailing}]]]$ ” is true in these worlds. Consequently, M is a model of the above program.

Two properties often associated to knowledge are (see the work of Fagin et al,²⁴) the *knowledge generalisation rule* and the *truth axiom of knowledge*. The knowledge generalisation rule states that knowledge present in all worlds is knowledge of all contexts. Formally:

$$\text{if } M \models \varphi \text{ then } M \models d[\varphi]$$

where $M \models \varphi$ means that φ is true in all worlds (φ is valid). If φ is valid, then d knows it to be true in all worlds. Let M be a model of the program “ $\text{doc}[\text{sailing sec}[\text{boats}]]$ ” (i.e. $M \models \text{this}[\text{doc}[\text{sailing sec}[\text{boats}]]]$). From Definition 2,

we have $M \models \text{doc}[\text{sailing}]$. Using the knowledge generalisation rule, we could derive that $M \models \text{sec}[\text{doc}[\text{sailing}]]$, which does not reflect the logical structure of the program. To overcome this problem, we introduce the notion of *context-validity* for structured context clauses:

$$M \models_C d[s[\varphi]] \iff M \models d[s[\varphi]] \text{ and } s \in \text{sub}(d)$$

A structured context is *context-valid* iff the structured context clause is valid and the logical structure of the program is reflected. The context-valid programs constitute a subset of the valid programs. The Definition 2 is extended as follows.

Definition 3 (Interpretation function \models - Extended Definition): *An interpretation structure M is a model of a program iff all context clauses of the program are valid and all structured context clauses are context-valid.*

In the remainder of this paper, we use \models to refer to a context-valid interpretation function.

The truth axiom of knowledge states that a context can only “know” what is true. Formally:

$$\text{if } (M, w) \models d[\varphi] \text{ then } (M, w) \models \varphi$$

Consider the program “doc[sailing]” and M as a model of this program. We have then from Definition 1 a set of possible worlds $W = \{w_{slot}, w_{this}, w_{doc}\}$ and the accessibility relations $R_{this} = \{(w_{slot}, w_{this})\}$ and $R_{doc} = \{(w_{this}, w_{doc})\}$, and truth value assignments $\pi(w_{doc})(sailing) = true$, $\pi(w_{slot})(sailing) = unknown$ and $\pi(w_{this})(sailing) = unknown$. It can be shown that by applying Definition 2 $(M, w_{this}) \models \text{doc}[\text{sailing}]$ although $(M, w_{this}) \not\models sailing$. Our knowledge modal operators do not follow the truth axiom of knowledge.³

This concludes our formalism (its syntax and semantics based on a Kripke structure) for representing the content and structural description of the objects forming a structured document. Objects are viewed as agents, referred to as contexts in this paper, that possess knowledge about their own content and their structure. The formalism is next extended so that the knowledge of an object is augmented with that of its structurally related objects.

3. Knowledge augmentation using modal operators

The semantics of programs as defined in the previous section formalises the content and structural knowledge of objects. The present section describes the formalisation of the knowledge augmentation process, necessary to allow for the retrieval of objects at varying granularity.

³The truth axiom of knowledge is satisfied if the accessibility relations are reflexive. Our accessibility relations are not reflexive from the way they are constructed (see Definition 1).

Consider the program “doc[sec1[sailing] sec2[boats]]”. A model for the program is based on an interpretation structure M defined upon a set of possible worlds $W = \{w_{slot}, w_{this}, w_{doc}, w_{sec1}, w_{sec2}\}$, the accessibility relations $R_{this} = \{(w_{slot}, w_{this})\}$, $R_{doc} = \{(w_{this}, w_{doc})\}$, $R_{sec1} = \{(w_{doc}, w_{sec1})\}$ and $R_{sec2} = \{(w_{doc}, w_{sec2})\}$, and the truth value assignments $\pi(w_{sec1})(sailing) = \pi(w_{sec2})(boats) = true$, and *unknown* in all other cases. Definition 2 yields $(M, w_{w_{doc}}) \models sec1[sailing]$ and $(M, w_{w_{doc}}) \models sec2[boats]$ characterising the knowledge content of sec1 and sec2 (sec1 knows sailing and sec2 knows boats), and $(M, w_{this}) \models doc[sec1[sailing]]$ and $(M, w_{this}) \models doc[sec2[boats]]$ characterising the structural knowledge of doc (doc knows that sec1 knows sailing and sec2 knows boats). Contexts sec1 and sec2 would be relevant to queries about “sailing” and “boats”, respectively, but none of the contexts sec1 and sec2 would be relevant to a query about “sailing and boats”. The supercontext doc should be relevant to such a query because it knows that one of its subcontext contains sailing and the other contains boats. This could be inferred if the knowledge content of doc is augmented with that of its subcontexts.

However, a knowledge augmentation process can lead to inconsistent knowledge. Consider the program “doc[sec1[sailing] sec2[not sailing boats]]”. Subcontext sec1 knows sailing; subcontext sec2 knows the opposite. The example makes evident that augmenting the content knowledge of doc by that of sec1 and sec2 leads to inconsistent knowledge regarding the proposition “sailing”, since we have evidence for true from sec1 and false from sec2. In the augmented context doc(sec1,sec2), “sailing” is therefore inconsistent.

To distinguish between the content knowledge of a context and its *augmented content knowledge*, we introduce the terminology of *augmented context* as opposed to *basic context*. Basic contexts are context names (e.g. doc, sec1 and sec2), whereas *augmented contexts* consist of a supercontext name and a list (group) composed of augmented contexts or basic contexts (e.g. doc(sec1(sec11),sec2)). A context clause with only basic contexts is called a *basic context clause* (e.g. “doc[sailing]” is a basic atomic context clause and “doc[sec[sailing]]” is a basic structured context clause). A context clause with augmented contexts is called an *augmented context clause* (e.g. “doc(sec1,sec2)[sailing]”).

The knowledge augmentation process *combines* knowledge of a supercontext with that of its subcontexts. Modal operators have been defined to formalise specific combination of knowledge; for example, common knowledge and distributed knowledge.²⁴ However, as discussed in previous work,¹⁰ these operators are not appropriate for a combination of knowledge that arises from a knowledge augmentation process. We therefore define other modal operators.

The first operator is the *united knowledge modal operator* denoted U_G , which is used to represent the combined knowledge of a group of contexts $G = (s1, \dots, sn)$ referred to as a *united context*. Here, we are aiming at capturing the combined knowledge of a group of context, not the knowledge of an augmented context.

Definition 4 (United knowledge operator U_G): Let $M = (W, \pi, \mathcal{R})$ be an interpretation structure defined upon a set Φ of propositions and a set \mathcal{C} of context names. Let w be a world in W , φ a proposition in Φ , and G a united context from \mathcal{C} .

$$\begin{aligned}
(M, w) \models U_G[1 \varphi] & : \iff \exists s \in G : (M, w) \models s[1 \varphi] \text{ and} \\
& \forall s \in G : ((M, w) \models s[1 \varphi] \text{ or } (M, w) \models s[0/0/0/1 \varphi]) \\
(M, w) \models U_G[0/1 \varphi] & : \iff \exists s \in G : (M, w) \models s[0/1 \varphi] \text{ and} \\
& \forall s \in G : ((M, w) \models s[0/1 \varphi] \text{ or } (M, w) \models s[0/0/0/1 \varphi]) \\
(M, w) \models U_G[0/0/1 \varphi] & : \iff (\exists s \in G : (M, w) \models s[1 \varphi] \text{ and} \\
& \exists s \in G : (M, w) \models s[0/1 \varphi]) \text{ or} \\
& \exists s \in G : (M, w) \models s[0/0/1 \varphi] \\
(M, w) \models U_G[0/0/0/1 \varphi] & : \iff \forall s \in G : (M, w) \models s[0/0/0/1 \varphi]
\end{aligned}$$

The proposition φ is true in a united context G iff there exists a context s in G that knows φ to be true and all contexts know φ to be true or unknown. The definition for false is analogous. The proposition φ is inconsistent iff there exist contexts that know φ to be true and false, or there exists a context in G that knows φ to be inconsistent. The proposition φ is unknown iff all contexts in G know φ to be unknown. For our example, we would obtain $(M, w_{doc}) \models U_{(\text{sec1}, \text{sec2})} [\text{boats}]$, $(M, w_{doc}) \models U_{(\text{sec1}, \text{sec2})} [\text{sailing}]$.

We define now the knowledge of augmented contexts through the introduction of an *augmented knowledge modal operator* A_{dG} where dG is an augmented context (d is the supercontext and G is the united context formed with the subcontexts of d , i.e. $G = \{s_i : s_i \in \text{sub}(d)\}$). The knowledge of the united context is combined with the knowledge of the supercontext. The augmented knowledge modal operator is therefore defined upon the united knowledge modal operator.

Definition 5 (Augmented knowledge operator A_{dG}): Let $M = (W, \pi, \mathcal{R})$ be an interpretation structure defined upon a set Φ of propositions and a set \mathcal{C} of context names. Let w be a world in W , φ a proposition in Φ , d, s_1, \dots, s_n contexts in \mathcal{C} such that $\text{sub}(d) = \{s_1, \dots, s_n\}$, and R_d the accessibility relation in \mathcal{R} associated with the supercontext d .

$$\begin{aligned}
(M, w) \models A_{d(s_1, \dots, s_n)}[1 \varphi] & : \iff \\
& ((M, w) \models d[1 \varphi] \text{ and} \\
& \forall w' \in R_d(w) : ((M, w') \models U_{(s_1, \dots, s_n)}[1 \varphi] \text{ or } (M, w') \models U_{(s_1, \dots, s_n)}[0/0/0/1 \varphi])) \text{ or} \\
& ((M, w) \models d[0/0/0/1 \varphi] \text{ and} \\
& \forall w' \in R_d(w) : (M, w') \models U_{(s_1, \dots, s_n)}[1 \varphi])
\end{aligned}$$

$$\begin{aligned}
(M, w) \models A_{d(s_1, \dots, s_n)}[0/1 \varphi] &: \iff \\
((M, w) \models d[0/1 \varphi] \text{ and} \\
\forall w' \in R_d(w) : ((M, w') \models U_{(s_1, \dots, s_n)}[0/1 \varphi] \text{ or } (M, w') \models U_{(s_1, \dots, s_n)}[0/0/0/1 \varphi])) \text{ or} \\
((M, w) \models d[0/0/0/1 \varphi] \text{ and} \\
\forall w' \in R_d(w) : (M, w') \models U_{(s_1, \dots, s_n)}[0/1 \varphi]) \\
(M, w) \models A_{d(s_1, \dots, s_n)}[0/0/1 \varphi] &: \iff \\
((M, w) \models d[0/0/1 \varphi] \text{ or} \\
\forall w' \in R_d(w) : (M, w') \models U_{(s_1, \dots, s_n)}[0/0/1 \varphi]) \text{ or} \\
((M, w) \models d[1 \varphi] \text{ and} \\
\forall w' \in R_d(w) : (M, w') \models U_{(s_1, \dots, s_n)}[0/1 \varphi]) \text{ or} \\
((M, w) \models d[0/1 \varphi] \text{ and} \\
\forall w' \in R_d(w) : (M, w') \models U_{(s_1, \dots, s_n)}[1 \varphi]) \\
(M, w) \models A_{d(s_1, \dots, s_n)}[0/0/0/1 \varphi] &: \iff \\
(M, w) \models d[0/0/0/1 \varphi] \text{ and} \\
\forall w' \in R_d(w) : (M, w') \models U_{(s_1, \dots, s_n)}[0/0/0/1 \varphi]
\end{aligned}$$

The proposition φ is true in an augmented context $d(s_1, \dots, s_n)$ iff (i) φ is true in the supercontext d and the united context (s_1, \dots, s_n) gives evidence for true or unknown, (ii) φ is unknown in the supercontext and the united context gives evidence for true. The definition for false is analogous. The proposition φ is inconsistent in $d(s_1, \dots, s_n)$ iff (i) φ is inconsistent in the supercontext or the united context gives evidence for inconsistent, (ii) φ is true (false) in the supercontext and the united context gives evidence for false (true). The proposition φ is unknown in $d(s_1, \dots, s_n)$ iff φ is unknown in the supercontext and the united context gives evidence for unknown.

Applying Definition 5 to our program “doc[sec1[sailing] sec2[boats]]”, we obtain $(M, w_{this}) \models A_{\text{doc}(\text{sec1}, \text{sec2})} [\text{sailing}]$ and $(M, w_{this}) \models A_{\text{doc}(\text{sec1}, \text{sec2})} [\text{boats}]$ because $(M, w_{doc}) \models U_{(s_1, s_2)} [\text{sailing}]$, $(M, w_{doc}) \models U_{(s_1, s_2)} [\text{boats}]$. By augmenting the content knowledge of doc by that of its subcontexts sec1 and sec2 , we arrive at a representation of the object represented by the context name doc that can be assessed relevant to a query about “sailing and boats”.

The definitions of the united and augmented knowledge modal operators allow the propagation of inconsistent knowledge. Inconsistency in a subcontext of a group is propagated to the united knowledge of the group, and inconsistency in the supercontext or the group of an augmented context leads to inconsistent knowledge in the augmented context. For example, a proceedings (the supercontext) may contain two documents (the subcontexts) that contradict each other. The augmented

knowledge of the supercontext is inconsistent. Combining this proceedings with another (e.g. to form a digital library) preserves this inconsistency.

4. Knowledge augmentation using truth value assignment functions

In the previous section, the knowledge augmentation process was defined through the introduction of modal operators. We defined two modal operators for representing united knowledge and augmented knowledge. The semantics of the knowledge modal operators is based on an interpretation structure defined upon a set of possible worlds, truth value assignment functions and accessibility relations. In this section, we define the knowledge augmentation process based on a graphical representation of worlds and accessibility relations formalised through the notions of G-world-trees. This follows the approach adopted by Fagin et al, ²⁴ where combination of knowledge was defined upon the notion of “G-reachability”. Two truth value assignment functions are defined upon G-world trees to characterise united knowledge and augmented knowledge.

The logical structure of a non-atomic object can be described as trees, referred to as *G-world-trees*. These tree structures are also reflected in the way possible worlds and accessibility relations are defined with respect to contexts.

An empty united context is denoted “()”, and for a context name s in \mathcal{C} , s is treated as augmented context with an empty united context (e.g. s is the same as $s()$).

Definition 6 (G-world-trees): *Let $M = (W, \pi, \mathcal{R})$ be an interpretation structure defined upon a set Φ of propositions and a set \mathcal{C} of context names. Let d be a context in \mathcal{C} with accessibility relation R_d in \mathcal{R} . Let w and w_0 be worlds in W . Let G_n be the united context (s_1, \dots, s_n) and let G_{n+1} be the united context (s_1, \dots, s_{n+1}) , for s_1, \dots, s_{n+1} in \mathcal{C} . The set of G-world-trees associated with a united context is defined inductively as follows:*

$$\begin{aligned} trees(w, G_{n+1}) := \{ & (w, S) \mid \exists S_n, t : (w, S_n) \in trees(w, G_n) \wedge \\ & t \in trees(w, s_{n+1}) \wedge S = S_n \cup \{t\} \} \end{aligned}$$

A tuple (w, S) is a G-world-tree of a world w and a united context G_{n+1} if there exists a set S_n such that (w, S_n) is a G-world-tree of the world w and the united context G_n , there exists a G-world-tree t of the world w and the context s_{n+1} , and $S = S_n \cup \{t\}$.

The G world-tree of a world w and an empty united context $()$ is $(w, \{\})$.

The set of trees associated with an augmented context is defined inductively as:

$$trees(w_0, dG) := \{ (w, S) \mid w \in R_d(w_0) \wedge (w, S) \in trees(w, G) \}$$

A tuple (w, S) is a G-world-tree of a world w_0 and an augmented context dG if $w \in R_d(w_0)$ and (w, S) is a G-world-tree of the world w and the united context G .

Using our example program “doc[sec1[sailing] sec2[boats]]” and the interpretation structure obtained in Section 3, we obtain the following G-world trees:

$$\begin{aligned}
trees(w_{sec1}, ()) &= \{(w_{sec1}, \{\})\} \\
trees(w_{doc}, sec1) &= \{(w, S) | w \in R_{sec1}(w_{doc}) \wedge (w, S) \in trees(w, ())\} \\
&= \{(w_{sec1}, \{\})\} \\
trees(w_{doc}, (sec1)) &= \{(w_{doc}, S) | \exists S_n, t : (w_{doc}, S_n) \in trees(w_{doc}, ()) \wedge \\
&\quad t \in trees(w_{doc}, sec1) \wedge S = S_n \cup \{t\}\} \\
&= \{(w_{doc}, \{(w_{sec1}, \{\})\})\} \\
trees(w_{doc}, (sec1, sec2)) &= \{(w_{doc}, S) | \exists S_n, t : (w_{doc}, S_n) \in trees(w_{doc}, (sec1)) \wedge \\
&\quad t \in trees(w_{doc}, sec2) \wedge S = S_n \cup \{t\}\} \\
&= \{(w_{doc}, (w_{sec1}, \{\}), (w_{sec2}, \{\}))\} \\
trees(w_{this}, doc(sec1, sec2)) &= \{(w, S) | w \in R_{doc}(w_{this}) \wedge (w, S) \in trees(w, (sec1, sec2))\} \\
&= \{(w_{doc}, (w_{sec1}, \{\}), (w_{sec2}, \{\}))\}
\end{aligned}$$

The G-world tree of an augmented context $d(s_1, \dots, s_n)$ formalises the accessibility of possible worlds associated with the subcontexts s_i from the possible world associated with the context d , and for each s_i , the accessibility of the possible worlds associated with the subcontexts of the s_i from the possible world associated with the context s_i , etc. This reflects the logical structure among d , its subcontexts s_i , the subcontexts of s_i , etc. We call \mathcal{T} the set of G-world trees associated with an interpretation structure M .

The next step is to define truth value assignment functions with respect to G-world trees to capture the logical structure among the contexts. First we define the truth value assignment function, referred to as *united truth value assignment function* associated with a united context to model the united knowledge of the united context. Here, the definition of the four truth values as sets (see Definition 1) is exploited to arrive at the *united truth value*.

Definition 7 (United truth value assignment function π_U): *Let $M = (W, \pi, \mathcal{R})$ be an interpretation structure defined upon a set Φ of propositions and a set \mathcal{C} of context names. Let \mathcal{T} the set of G-world trees associated with M . Let (w, S) be a G-world-tree in \mathcal{T} of world $w \in W$ and a united context G from \mathcal{C} . The united truth value function $\pi_U : \mathcal{T} \rightarrow (\oplus \rightarrow \{\text{true}, \text{false}, \text{inconsistent}, \text{unknown}\})$ of a G-world-tree (w, S) is defined as the union of the truth value functions $\pi(w')$ where the worlds w' are the roots of the subtrees $(w', \{\})$ in the set S . Formally, for all proposition φ in Φ :*

$$\pi_U((w, S))(\varphi) := \bigcup_{(w', \{\}) \in S} \pi(w')(\varphi)$$

For an empty set S , the united truth value $\pi_U((w, \{\}))(\varphi) := \text{unknown}$.

Consider the united truth value of the united context (sec1,sec2). Given the accessibility relations $R_{sec1} = \{(w_{doc}, w_{sec1})\}$ and $R_{sec2} = \{(w_{doc}, w_{sec2})\}$, we obtain for instance⁴:

$$\begin{aligned} \pi_U((w_{doc}, \{(w_{sec1}, \{\}), (w_{sec2}, \{\})\}))(sailing) &= \\ \pi(w_{sec1})(sailing) \cup \pi(w_{sec2})(sailing) &= \\ true \cup unknown &= \\ true \end{aligned}$$

This means that with respect to the logical structure between sec1 and sec2, the united truth value of sailing in w_{doc} is *true*.

We define now the truth value assignment function, referred to as *augmented truth value assignment function*, associated with an augmented context to model its augmented knowledge. Analogously to united truth value, the function is defined upon G-world-trees thus capturing the logical structure between the contexts, and the definition of the four truth values as sets is again exploited to arrive at the *augmented truth value*.

Definition 8 (Augmented truth value assignment function π_A): *Let $M = (W, \pi, \mathcal{R})$ be an interpretation structure defined upon a set Φ of propositions and a set \mathcal{C} of context names. Let \mathcal{T} the set of G-world trees associated with M . Let d be a context in \mathcal{C} and G a united context in \mathcal{C} . Let (w, S) be a G-world-tree in \mathcal{T} of world w_0 in W and the augmented context dG . The augmented truth value function $\pi_A : \mathcal{T} \rightarrow (\oplus \rightarrow \{\text{true}, \text{false}, \text{inconsistent}, \text{unknown}\})$ of a G-world-tree (w, S) is defined as the union of the truth value function $\pi(w)$ of world w and the united truth value function $\pi_U((w, S))$ of the G-world-tree (w, S) . Formally, for all φ in Φ :*

$$\pi_A((w, S))(\varphi) := \pi(w)(\varphi) \cup \pi_U((w, S))(\varphi)$$

Consider the G-world-tree $S = (w_{doc}, \{(w_{sec1}, \{\}), (w_{sec2}, \{\})\})$ for the world w_{this} and the augmented context $\text{doc}(\text{sec1}, \text{sec2})$. We obtain for instance

$$\begin{aligned} \pi_A((w_{doc}, S))(sailing) &= \\ \pi(w_{doc})(sailing) \cup \pi_U((w_{doc}, S))(sailing) &= \\ \pi(w_{doc})(sailing) \cup \pi(w_{sec1})(sailing) \cup \pi(w_{sec2})(sailing) &= \\ unknown \cup true \cup unknown &= \\ true \end{aligned}$$

⁴We recall that $true = \{t\}$ and $unknown = \{\}$ so $true \cup unknown = \{t\} \cup \{\} = \{t\} = true$.

World w_{doc} is associated with the context doc . For the basic context doc , sailing is *unknown* in w_{doc} , whereas, for the augmented context $doc(sec1,sec2)$, sailing is *true* in w_{doc} , thus modelling augmented knowledge.

5. Equivalence of the two approaches

Sections 3 and 4 present two approaches that formalise the knowledge augmentation process. The two approaches are based on different formalisations of united knowledge (Definitions 4 and 7) and augmented knowledge (Definitions 5 and 8). In the present section, we formally show that the two approaches are equivalent, that is, they lead to the same knowledge augmentation process. First, we show that the definitions of united truth value assignment function (Definition 7) and united knowledge modal operator (Definition 4) lead to the same truth value assignment to a proposition in a united context. We only provide the proof for the truth value *true*.

Theorem 1 (United knowledge): *Let $M = (W, \pi, \mathcal{R})$ be an interpretation structure defined upon a set Φ of propositions and a set \mathcal{C} of context names. Let \mathcal{T} the set of G -world trees associated with M . In a world $w \in W$, a proposition $\varphi \in \Phi$ is true in a united context G iff for each tree (w, S) of world w and united context G , the united truth value equals true. Formally $\forall s$ in $G : R_s(w) \neq \{\}$:*

$$(M, w) \models U_G[\varphi] \iff \forall (w, S) \in \text{trees}(w, G) : \pi_U((w, S))(\varphi) = \text{true}$$

We prove the theorem by contradiction and induction over united contexts.

Proof: Let G be (s1), i.e. a united context with one context. Assume that there exists a tree (w, S) for the world w and the united context (s1) such that $\pi_U((w, S))(\varphi) \neq \text{true}$. According to Definition 7, $\bigcup_{(w', \{\}) \in S} \pi(w')(\varphi) \neq \text{true}$. According to Definition 6 (G-world-tree), for a united context with one context, we obtain:

$$\begin{aligned} \text{trees}(w, (s1)) &= \{(w, S) \mid \exists t : t \in \text{trees}(w, s1) \wedge S = \{t\}\} \\ \text{trees}(w, s1) &= \{(w', \{\}) \mid w' \in R_{s1}(w)\} \end{aligned}$$

The set S contains one tree $(w', \{\})$, $w' \in R_{s1}(w)$. Since $\bigcup_{(w', \{\}) \in S} \pi(w')(\varphi) \neq \text{true}$, there exists a world w' such that $w' \in R_{s1}(w)$ and the truth value $\pi(w')(\varphi) \neq \text{true}$. According to Definition 2, $s1[\varphi]$ is not true in world w . According to Definition 4, $U_{(s1)}[\varphi]$ is not true in world w .

We prove now the reverse. Assume that $U_{(s1)}[\varphi]$ is not true in world w . According to definition 4 (united knowledge), $s1[\varphi]$ is not true in world w . According to Definition 2 (interpretation function), there exists a world w' such that $w' \in R_{s1}(w)$ and the truth value $\pi(w')(\varphi) \neq \text{true}$. According to Definition 6 (G-world-tree), for a united context with one context, we obtain:

$$\text{trees}(w, (s1)) = \{(w, S) \mid \exists t : t \in \text{trees}(w, s1) \wedge S = \{t\}\}$$

$$\text{trees}(w, s_1) = \{(w', \{\}) \mid w' \in R_{s_1}(w)\}$$

The set S contains one tree $(w', \{\})$, $w' \in R_{s_1}(w)$. Since there exists a world w' such that $w' \in R_{s_1}(w)$ and the truth value $\pi(w')(\varphi) \neq \text{true}$, $\bigcup_{(w', \{\}) \in S} \pi(w')(\varphi) \neq \text{true}$. According to Definition 7 (united truth value), there exists a tree (w, S) for the world w and the united context (s_1) such that $\pi_U((w, S))(\varphi) \neq \text{true}$.

Assume that the theorem holds for a united context G_n with n contexts.

$$(M, w) \models U_{G_n}[\varphi] \iff \forall (w, S) \in \text{trees}(w, G_n) : \pi_U((w, S))(\varphi) = \text{true}$$

Adding one context $sn+1$ leads to the united context G_{n+1} , and we prove:

$$(M, w) \models U_{G_{n+1}}[\varphi] \iff \forall (w, S) \in \text{trees}(w, G_{n+1}) : \pi_U((w, S))(\varphi) = \text{true}$$

Assume that there exists a tree (w, S) for the world w and the united context G_{n+1} such that $\pi_U((w, S))(\varphi) \neq \text{true}$. From Definition 6 (G-world-tree), we obtain the G-world-trees of the united context G_n and context $sn+1$.

$$\begin{aligned} (w, S) \in \text{trees}(w, G_{n+1}) &\iff \\ (w, S_n) \in \text{trees}(w, G_n) \wedge t \in \text{trees}(w, sn+1) \wedge S &= S_n \cup \{t\} \end{aligned}$$

According to Definition 7, $(\pi_U((w, S_n)) \cup \pi_U((w, \{t\}))) (\varphi) \neq \text{true}$ since the united truth value $\pi_U((w, S))(\varphi) \neq \text{true}$. From Definition 4 (united knowledge), we obtain that for $U_{G_n}[\varphi]$ being true, φ is neither true nor unknown in $(sn+1)$, and for φ being unknown in G_n , φ is not true in $(sn+1)$. Thus, $U_{G_{n+1}}[\varphi]$ is not true in world w .

Assume that $U_{G_{n+1}}[\varphi]$ is not true in world w . From Definition 4 (united knowledge), we obtain that for $U_{G_n}[\varphi]$ being true, φ is neither true nor unknown in $(sn+1)$, and for φ being unknown in G_n , φ is not true in $(sn+1)$. From Definition 6 (G-world-tree), we obtain the G-world-trees of the united context G_n and context $sn+1$.

$$\begin{aligned} (w, S) \in \text{trees}(w, G_{n+1}) &\iff \\ (w, S_n) \in \text{trees}(w, G_n) \wedge t \in \text{trees}(w, sn+1) \wedge S &= S_n \cup \{t\} \end{aligned}$$

According to Definition 7, the united truth value $\pi_U((w, S))(\varphi) \neq \text{true}$, since $(\pi_U((w, S_n)) \cup \pi_U((w, \{t\}))) (\varphi) \neq \text{true}$. Thus, there exists a tree (w, S) for the world w and the united context G_{n+1} such that $\pi_U((w, S))(\varphi) \neq \text{true}$. \square

We have shown the equivalence between the united knowledge modal operator (Definition 4) and the united truth value assignment function (Definition 7). We turn now to augmented knowledge, where we show that the definitions of the augmented truth value assignment (Definition 8) and the augmented knowledge modal operator (Definition 5) lead to the same truth value of a proposition in an augmented context. As for united knowledge, we prove the equivalence for the truth value *true* only.

Theorem 2 (Augmented knowledge): *Let $M = (W, \pi, \mathcal{R})$ be an interpretation structure defined upon a set Φ of propositions and a set \mathcal{C} of context names. Let*

\mathcal{T} the set of G -world trees associated with M . In a world w_0 in W , a proposition φ in Φ is true in an augmented context dG iff for each tree (w, S) of world w_0 and an augmented context dG , the augmented truth value equals true.

$$(M, w_0) \models dG[\varphi] \iff \forall (w, S) \in \text{trees}(w_0, dG) : \pi_A((w, S))(\varphi) = \text{true}$$

We prove the theorem by contradiction for an augmented context dG .

Proof: Assume that there exists a tree $(w, S) \in \text{trees}(w_0, dG)$ for the world w_0 and the augmented context dG such $\pi_A((w, S))(\varphi) \neq \text{true}$. According to Definition 8 (augmented truth value), $(\pi(w) \cup \pi_U((w, S))) (\varphi) \neq \text{true}$. It follows that if $\pi(w)(\varphi) = \text{true}$, then $\pi_U((w, S))(\varphi)$ is neither *true* nor *unknown*, and if $\pi(w)(\varphi) = \text{unknown}$, then $\pi_U((w, S))(\varphi) \neq \text{true}$. The same applies analogously for switched roles of $\pi(w)$ and $\pi_U((w, S))$. It follows that if φ is true in d for world w_0 ($(M, w_0) \models d[1 \varphi]$), then there exists $w \in R_d(w_0)$ such that φ is neither *true* nor *unknown* in the united context G for world w , and if φ is *unknown* in d for world w_0 ($(M, w_0) \models d[0/0/0/1\varphi]$), then there exists $w \in R_d(w_0)$ such that $U_G[\varphi]$ is not true in world w . The same applies analogously for switched roles of $d[\varphi]$ and $U_G[\varphi]$. Thus, $dG[\varphi]$ is not true in world w_0 .

Assume that $dG[\varphi]$ is not true in world w_0 . According to Definition 5 (augmented knowledge), if φ is true in d for world w_0 ($(M, w_0) \models d[1 \varphi]$), then there exists $w \in R_d(w_0)$ such that φ is neither *true* nor *unknown* in the united context G for world w , and if φ is *unknown* in d for world w_0 ($(M, w_0) \models d[0/0/0/1\varphi]$), then there exists $w \in R_d(w_0)$ such that $U_G[\varphi]$ is not true in world w . The same applies analogously for switched roles of $d[\varphi]$ and $U_G[\varphi]$. It follows that if for all $w \in R_d(w_0)$, $\pi(w)(\varphi) = \text{true}$, then there exists a tree (w, S) for the united context G such that $w \in R_d(w_0)$ and $\pi_U((w, S))(\varphi)$ is neither *true* nor *unknown*, and if for all $w \in R_d(w_0)$, $\pi(w)(\varphi) = \text{unknown}$, then there exists a tree (w, S) for the united context G such that $\pi_U((w, S))(\varphi) \neq \text{true}$. The same applies analogously for switched roles of $\pi(w)$ and $\pi_U((w, S))$. It follows that there exists a world $w \in R_d(w_0)$ such that $(\pi(w) \cup \pi_U((w, S))) (\varphi) \neq \text{true}$. Thus, there exists a tree $(w, S) \in \text{trees}(w_0, dG)$ for the world w_0 and the augmented context dG such $\pi_A((w, S))(\varphi) \neq \text{true}$. \square

The section shows, via (a sample of) its proofs, that our first formalism (Section 3) is compatible with the well-known framework defined by Fagin et al.,²⁴ for modelling the combination of knowledge (Section 4).

6. Conclusion and future work

This paper describes a formal model for representing the content and logical structure of structured documents for the purpose of their retrieval. To obtain a representation that allows for the retrieval of objects of varying granularity, the content of an object is viewed as the knowledge contained in that object, and the structure among objects is captured by a process of knowledge augmentation, which leads to a representation of the object that is based on its own content and that of its

structurally related objects.

The model is based on a Kripke structure in which the structure of a document is reflected. The objects composing the document are viewed as contexts (agents) that possess knowledge about their own content and their structure. This knowledge is augmented with that of its structurally related objects. The model is based on the definitions of knowledge modal operators, augmented knowledge operators and augmented truth value assignment functions, formalised upon an interpretation structure (possible worlds, truth value assignments and accessibility relations). The model uses four truth values to capture incompleteness and inconsistency. The knowledge augmentation process is formalised through the use of modal operators to represent the content knowledge of an object augmented with that of its structurally related objects.

We also formalise the knowledge augmentation process using the framework defined by Fagin et al,²⁴ to model the combination of knowledge. With this approach, the logical structure of an object is represented through the use of G-world-trees upon which augmented truth values are defined.

The formalism described in this paper is the basis of the development of a model for structured document retrieval that allows for the retrieval of objects of varying granularity. Future work will present the representation of the uncertainty inherent to the information retrieval process,^{27,28} which will be used to estimate the degree to which an object is relevant to a query.

The knowledge augmentation process has been implemented in the HySpirit platform, an information retrieval experimental platform available at Queen Mary, which allows describing, implementing and evaluating indexing and retrieval strategies. The implemented knowledge augmentation process has been applied in two projects at Queen Mary⁵: (1) the focused retrieval of structured documents, which aims at returning best entry points in a structured document; and (2) the retrieval of video data annotated with MPEG-7, the metadata standard for describing the content and structure of multimedia data. In both cases, the knowledge augmentation process provided a means to effectively retrieve objects of varying granularity.

In future work, we will compare our model to other models for structured documents. This comparison will be performed within the INEX⁶ initiative, which aims at providing a common forum to evaluate approaches for XML retrieval using a large collection of XML documents.

References

1. G. Bordogna and G. Pasi. Flexible querying of structured documents. In *Proc. of Flexible Query Answering Systems (FQAS)*, pages 350–361, Warsaw, Poland, 2000.
2. Y. Chiamella. Browsing and querying: two complementary approaches for multimedia information retrieval. In *Proc. of Hypermedia - Information Retrieval - Multimedia*, Dortmund, Germany, 1997.

⁵see <http://qmir.dcs.qmul.ac.uk>.

⁶<http://qmir.dcs.qmul.ac.uk/INEX>

3. Y. Chiamarella, P. Mulhem, and F. Fourel. A model for multimedia information retrieval. Technical Report Fermi ESPRIT BRA 8134, University of Glasgow, 1996.
4. T.P. Chinenyanga and N. Kushmerick. Expressive retrieval from XML documents. In *Proc. of Research and Development in Information Retrieval*, pages 163–171, New Orleans, USA, 2001.
5. N. Fuhr and K. Grossjohann. XIRQL: A query language for information retrieval in XML documents. In *Proc. of Research and Development in Information Retrieval*, pages 172–180, New Orleans, USA, 2001.
6. M Gery and J-P Chevallet. Toward a structured information retrieval system on the web: Automatic structure extraction of web pages. In *International Workshop on Web Dynamics*, London, 2001.
7. E. Kotsakis. Structured information retrieval in XML documents. In *Proc. of Symposium on Applied Computing (SAC02)*, Madrid, Spain, 2002.
8. S. Abiteboul, S. Cluet, V. Christophides, T. Milo, G. Moerkotte, and J. Simeon. Querying documents in object databases. *Int. J. on Digital Libraries*, 1:1–9, 1997.
9. T. Rölleke and N. Fuhr. Retrieval of complex objects using a four-valued logic. In *Proc. of Research and Development in Information Retrieval*, pages 206–214, Zurich, Switzerland, 1996.
10. T. Rölleke. *POOL: Probabilistic Object-Oriented Logical Representation and Retrieval of Complex Objects - A Model for Hypermedia Retrieval*. PhD thesis, University of Dortmund, Germany, 1999.
11. J. Callan. Passage-level evidence in document retrieval. In *Proc. of Research and Development in Information Retrieval*, pages 302–310, Dublin, Ireland, 1994.
12. R. Wilkinson. Effective retrieval of structured documents. In *Proc. of Research and Development in Information Retrieval*, pages 311–317, Dublin, Ireland, 1994.
13. G. Salton, J. Allan, and C. Buckley. Approaches to passage retrieval in full text information systems. In *Proc. of Research and Development in Information Retrieval*, pages 49–58, Pittsburgh, USA, 1993.
14. I.A. Macleod. Storage and retrieval of structured documents. *Information Processing and Management*, 26(2):197–208, 1990.
15. F.J. Burkowski. Retrieval activities in a database consisting of heterogeneous collections of structured texts. In *Proc. of Research and Development in Information Retrieval*, pages 112–125, Copenhagen, Denmark, 1992.
16. G. Navarro and R. Baeza-Yates. A language for queries on structured and contents of textual databases. In *Proc. of Research and Development in Information Retrieval*, pages 93–101, Seattle, USA, 1995.
17. M. Lalmas. Dempster-Shafer's theory of evidence applied to structured documents: modelling uncertainty. In *Proc. of Research and Development in Information Retrieval*, pages 110–118, Philadelphia, PA, USA, 1997.
18. C. Baumgarten. A probabilistic model for distributed information retrieval. In *Proc. of Research and Development in Information Retrieval*, pages 258–266, Philadelphia, USA, 1997.
19. S.H. Myaeng, D. H. Jang, M. S. Kim, and Z. C. Zhoo. A flexible model for retrieval of SGML documents. In *Proc. of Research and Development in Information Retrieval*, pages 138–145, Melbourne, Australia, 1998.
20. M.E. Frisse. Searching for information in a hypertext medical handbook. *Communications of the ACM*, 31(7):880–886, 1988.
21. N. Fuhr, N Goevert, and T Rölleke. Dolores: A system for logic-based retrieval of multimedia objects. In *Proc. of Research and Development in Information Retrieval*, Australia, 1999.
22. M. Lalmas and E. Moutogianni. A Dempster-Shafer indexing for the focussed retrieval

- of hierarchically structured documents: Implementation and experiments on a web museum collection. In *RIAO*, Paris, France, 2000.
23. B.F. Chellas. *Modal Logic*. Cambridge University Press, 1980.
 24. R. Fagin, J. Harpen, J. Moses, and M. Vardi. *Reasoning about Knowledge*. MIT Press, Cambridge, Massachusetts, 1995.
 25. N. Belnap. A useful four-valued logic. In J. Dunn and G. Epstein, editors, *Modern Uses of Multiple-valued Logic*. Reidel, Dordrecht, 1977.
 26. M. Lalmas, T. Rölleke, and N. Fuhr. Intelligent retrieval of hypermedia documents. Chapter in *Intelligent Exploration of the Web*, 2002. In Press.
 27. R. Baeza-Yates and B. Ribeiro-Neto. *Modern Information Retrieval*. Addison Wesley, 1999.
 28. N. Fuhr. Probabilistic models in information retrieval. *The Computer Journal, Special Issue*, 35(3):243–255, 1992.