# A FRAMEWORK FOR SEMANTICALLY RICH LEGAL DOCUMENTS AND APPLICATIONS

John Murphy and Robert Steele

*Department of Computer Science, University of Technology Sydney, PO Box 123 Sydney, 2007, Australia*

Keywords:     Semantic Web Application, Web Standards, Document Annotation, RDF, OWL.

Abstract:     This paper describes the use of Semantic Web technologies for enriching legal documents within a framework that provides a base for the development of applications to reduce problems experienced by legal practitioners in the administration of justice. The framework builds upon the terms and concepts found within structured information sources in the document centric legal domain. Contributions of this paper include (1) the description of a bottom-up ontology development approach utilizing the terms and structures within legal documents; (2) an explanation of how this legal knowledge can be captured from existing document resources and represented using Semantic Web languages; (3) an analysis of the comparative benefits of this bottom-up approach versus existing top-down approaches and (4) concrete examples of how to represent such legal knowledge using OWL/ RDF. We argue such bottom-up legal ontology development grounded in the terms and concepts found in existing legal resources will support direct and easier annotation of resources with the semantics of legal concepts and hence provide a base for the easier development of a legal applications layer.

## 1 INTRODUCTION

Recent advances in semantic web technologies have provided languages, tools and infrastructure that facilitate the application of these technologies to problems in existing domains. The maturing of languages such as OWL and RDF(S) with patterns of best practice has enabled the application of Semantic Web technologies to domains that have the characteristics of being document centric with distributed, heterogeneous and structured information sources. In this paper, we report on an application of Semantic Web technology to the legal domain and demonstrate the use of these technologies to enrich the content of legal documents and reduce the problems faced by legal practitioners dealing with complex issues using information sources spanning multiple documents.

The domain of law shares similar characteristics to the World Wide Web (WWW) because legal administration is primarily document driven, with many interrelated documents and independent authors. Legal document publishing is characterized by different types of related documents produced by different organizations that have a mix of semi-structured and free text parts. Unlike the WWW, the legal domain is underpinned by a formal legal vocabulary

that significantly reduces the problems of ambiguity. The characteristics of formal structure, reduced ambiguity and its distributed heterogeneous nature make the legal domain a strong candidate for the application of Semantic Web technologies and standards that encourage the integration of such information resources. Web technologies have evolved to overcome problems of distributed heterogeneous information sources by standardizing access to information resources. The Semantic Web technologies of RDF(S) and OWL have further enhanced the ability to process distributed information resources by enabling the standardization of the description of a documents structure and the meaning of its content.

Semantic Web technologies enable a standardized and richer description of the meaning of legal documents by annotation of the documents terms in relation to an ontology. Traditional methods of describing the meaning of documents are achieved through external categorization systems such as indexes and digests, that assign meaning to the whole document using tree like structures limited to expressing the single *is_kind_Of* relation between parent and child concepts. OWL and RDF(S) are languages that enable the formal definition of legal concepts and the relationships between those concepts using logical constructs

as well as providing the means to directly annotate those concepts within a document. Further benefits of OWL are obtained by reasoning using the logical constructs from Description Logic that makes explicit the meaning that is implicit in the document. An added benefit of using web standards with public legal documents is that it makes the information more accessible to the wider population and provides a platform that encourages broader participation in advancing research and application development in the legal publishing domain. Many courts and legislators are currently publishing directly on-line or sending their documents to commercial legal publishers or organizations such as AUSTLII (AustLII, 2007) for publication on the web.

Semantic Web technologies can be applied to solve current restrictions on searching and understanding information. Searching legal documents can be extended from "bag of words" or terms to searching by concepts. The development of multiple problem specific ontologies enables different perspectives of domain knowledge to be used by applications for searching and navigation of a set of related documents in a more focused and better targeted manner and provides the structure for understanding those documents from a problem specific point of view.

This paper demonstrates the use of Semantic Web technologies to deal with some of the problems faced by legal practitioners in the legal domain. It describes related semantic technology research in law, explains where our research sits with current legal ontology development and describes the approach we have taken to applying the Semantic Web technologies of OWL and RDF to the enrichment of legal documents.

## 2 BACKGROUND

**Ontology Development.** A significant part of legal semantic research has been devoted to top-down ontology development to produce re-usable shared ontologies. The current top-down approaches to legal ontology development have evolved from legal theories such as those of (Kelsen, 1991), (Hart, 1961) and (Bentham, 1970). These theories formed the basis of legal ontologies such as FOLaw (Valente, 1995) that represent functional explanations of "how the law works". The next evolution in top-down legal ontology development was the move toward representing "what the law is" by capturing commonsense concepts required for most applications of legal reasoning. LRI Core (Breuker and Hoekstra, 2004) is an example of a core commonsense legal ontology that

reflects the roles, objects, functions and actions of the legal system in society and are designed to facilitate reasoning and argumentation.

**Conceptual Search.** Applications of semantic technology to conceptual searching are based on the idea that established legal concepts can be modeled and used for document retrieval systems that locate information that conforms to a specific legal concept. Current searching techniques often use syntactic or statistical analysis of terms rather than concepts or semantic relations between terms. This often leads to document retrieval that is not ranked by the expected concept of relevance (van Noortwijk et al., 2005). For systems to achieve the goal of conceptual retrieval they require knowledge of legal concepts and issues as well as knowledge of the relationships of these concepts to cases and legislation (Hafner, 1987). Our framework establishes a base that applications can build upon to help attain this goal.

**Case Modeling.** Modeling the relations between legal principles is an important contribution to the task of legal research in its own right, because some of the associations between principals of law, legal issues and facts in a sub-domain of law are not always obvious to legal practitioners (Atkinson et al., 2005). Recent developments in the modeling of legal argumentation provide a structure for the development of ontologies to specifically capture this important aspect of judicial process.

**Legal Publishing on the Web.** Public legal documents such as court judgements and legislation are currently being published on websites such as AUSTLII (AustLII, 2007). AUSTLII provides free and open access to public legal documents on the web using an index based on word occurrences and databases of jurisdictions and parliaments for the publication of case law and legislation respectively. The AUSTLII case law databases include hypertext links to other cases, Acts of parliament and sections in Acts.

In contrast to other frameworks that extract concepts from the text using automated techniques (Biagioli et al., 2005) and (Casellas et al., 2006), our work concentrates on enriching existing legal documents and the relationships between and within those documents focusing on those parts that provide procedural context and summarization of the main legal points as well as explicit external links that provide context and supplementary information. Current approaches to legal publishing on the web are generally restrained to a document locator and optionally an anchor point

within the document. We argue there are significant benefits to be gained in the fields of searching, document retrieval, navigation and analysis of legal documents by enriching the documents structure, content and links by enhancing or extending them with deeper semantics. This we argue, will deliver more complex and multi-faceted documents for analysis and shift the development focus toward the concepts within documents rather than the document as a whole being the fundamental level of conceptualization.

# 3 A FRAMEWORK FOR SEMANTICALLY RICH LEGAL DOCUMENTS AND APPLICATIONS

The enrichment of legal documents through the application of Semantic Web technologies begins with the development of ontologies and the subsequent annotation of the existing legal documents. A bottom-up approach to development is used taking advantage of the structures and terminology found within legal documents to build higher level legal concepts that can be re-used in applications for document retrieval and navigation, analysis and domain understanding, document management, automated processing and legal reasoning.

A good example of a structured legal document is the common law *Case Report* that will be used throughout this paper to give examples of our ontology development approach, legal document annotation process and the types of applications supported by the framework. The *Case Report* document is an important contributor to traditional paper-based legal knowledge that contains the judgements of the court as well as more formally structured parts contributed by court reporters and editors as demonstrated in the *Case Report* extract presented in figure 1.



KR v KR
[2005] NSWCA 411
Court of Appeal: Hodgson JA, Bryson JA and Hunt A-JA
22 September, 25 November 2005

*Family Law --- Defacto relationships --- Adjustment of property interests --- Relevant Property --- Property (Relationships Act 1984, s 20.*

*Held*: (1) The exercise of jurisdiction under s 20 ...
　　　(2) The court may take a global view ...
　　　*AAA v BBB* (1998) 23 Fam LR 716, affirmed.
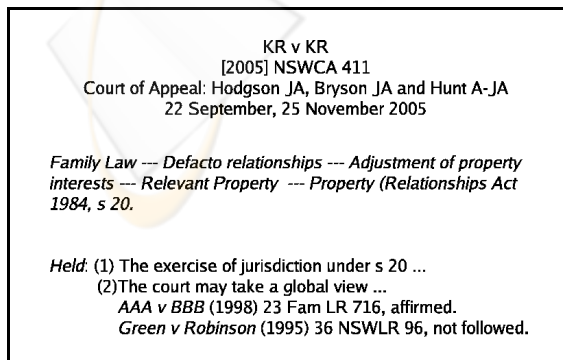　　　*Green v Robinson* (1995) 36 NSWLR 96, not followed.

Figure 1: Case Report Extract.

The extract is taken from the *Head Note* part of the *Case Report* and shows a structured mix of judicial case management information in the upper portion and information concerning legal principles in the portion marked "Held". The *Head Note* in addition to the court judgement, incorporates elements of judicial process such as hearing dates, parties and their advocates, the presiding judge, the type of court, descriptions of the nature of the case and cases cited. Furthermore the *Case Report* includes more complex parts such as the *"Catch Words"* block that is a hierarchical categorization of key terms derived from Legal Digests and a part that summarizes the important legal issues dealt with by the case. Mixed with the structured parts of the document are links to external unifying sources such as common legal phrases, terms from legal thesauri as well as links to relevant legislation.

The framework is composed of the three stages of 1) ontology development, 2) document annotation and 3) semantic application re-use. During the ontology development stage, higher level legal concepts are evolved from the explicit legal information found in the document text. The annotation stage involves the semi-automated annotation of the legal text using existing mark-up and knowledge of domain experts such as *Case Report* editors whose current role includes the identification and summarization of important legal principles. The semantic application stage involves the selection of relevant parts of the ontologies for re-use by specific applications for purposes such as document retrieval and management and legal reasoning.

## 3.1 Ontology Development

The framework adopts a bottom-up approach to ontology development using the structures and concepts contained within legal documents as the starting point. This approach to ontology development begins with the identification of concepts that are explicit in the legal documents including the names and roles of the parties, the case citation, date of hearing, the judges, the case and legislative citations, the advocates and references to the sections of any relevant legislation. Figure 2 shows the conceptual model of the *Case Report* document from which more abstract concepts of law are develop to:

- represent the administrative rules and processes of the court,

- encode the links between facts and legal principles,

- represent the roles of the participants in litigation,

- encode the links between the different types and level of courts,

- encode the links between legal principles, judgements and previous cases incorporating the concepts of *Following*, *Not Following*, *Affirming*, *Applying* and *Reversing* previous decisions and

- accommodate the representation of a sub-domain of law using a structure that facilitates understanding of the sub-domain. For example, the sub-domain of negligence within the laws of "Tort", would require a structure that encompasses the elements and sub-elements of negligence, "duty of care", "existence and breach of the duty of care" and "damage suffered as a result of the breach".
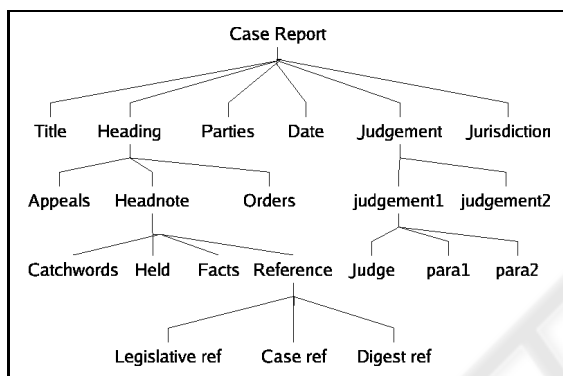


Figure 2: Case Report Concepts.

Bottom-up ontology development starts with the selection and definition of of the most specific concepts at the bottom of the hierarchy. The design advances toward more abstract concepts by grouping the lower concepts into more general concepts and by the application of rules to the lower concepts for the construction of higher level concepts. This bottom-up approach yields the benefits of accessibility and familiarity to the domain experts that annotate legal documents.

An example of developing higher level legal concepts from specific concepts explicit in the document is the *Considered Case* concept that is implied by the citation link identified in the *Case Report* extract in figure 1 as "AAA v BBB (1998) 23 FAM LR 716, affirmed". This is the citation of a case that the judge took into consideration when deciding the *KR v KR* case specifically *Affirming* the previous decision. The *Considered Case* is a higher level legal concept and judicial mechanism that supports the related concept of *Case Authority* in common law that encompasses both the adoption or rejection of legal precedent and defined by the (NYSUC, 2007) as a:

previously adjudged action or decision on

same or similar point, serving as a rule or example for present guidance.

*Case Authority* can be interpreted as the application of cumulative weightings of importance to a case over time using the decisions of subsequent judgements that have considered the case and *Followed*, *Not Followed*, *Reversed* or *Applied* the case rulings.

A number of benefits can be obtained using the framework that include:

- Accessibility by legal knowledge experts and legal practitioners because the approach builds upon the direct annotation of the original document source and the legal domain is a knowledge intensive domain that stores its knowledge in documents that can be marked-up using XML. This direct access to the source by means of a readily available technology gives non-technical personnel easier access to the source and enables broader participation in the capture of knowledge from the source documents.

- Familiarity of concepts used by domain experts because top-down approaches require knowledge of higher level concepts defined by external parties that are often highly abstract with subtle distinction that are not easy to understand and are difficult to use as a starting point for development.

- Better management of the development process by encouraging partitioning of concepts along familiar domain and task specific lines.

Our research project is developing ontologies to represent the knowledge of *Case Report* structures and case reporting rules with a mixture of complex and simple relations between legal concepts. The ontologies are designed to allow extension for the capture of complex legal principles and assist with tasks such as indexing, extraction and summarization of legal issues and the provision of an entry point for domain understanding and conceptual search for legal practitioners.

## 3.2 Annotation

After the required legal concepts are defined in ontologies, the lower level concepts are identified in the documents and annotated. This process can be semi-automated using existing mark-up to assist the identification of terms. Lower level ontology concepts can often be mapped directly to explicitly marked-up terms in the text whilst other mark-up that identifies document structures can be used to locate terms. Our approach to the ontology development yields the

benefit of starting with the most specific concepts that are explicit within the legal document or can be easily identified with human assistance.

A feature of legal documents that makes them more conducive to a bottom-up approach is that much of the publicly available legal documentation such as *Case Reports* and legislation already have their structure and many of the lower level explicit concepts marked-up in XML or SGML. Many jurisdictions have marked-up their court documents at the initial point of production by the the judge or judges associate and their work has been further expanded by *Case Reporters* who add value to the judgements by making explicit many concepts in the judgement document that are implicit. A benefit of the bottom-up approach is that it only requires the mark-up of the lower level legal concepts which facilitates the annotation of this large body of existing legal documentation through automatic or semi-automatic transformations of the existing XML or SGML.

Further benefits of direct annotation of legal documents are that it

- removes the impedance of external referencing of sources found in other types of knowledge stores such as databases or technologies acting external to the source such as Xpointer.

- relieves the complications associated with automatic identification of concepts using Natural Language Processing by relying on the intervention of human legal experts such as drafters of legislation and *Case Report* editors.

A final benefit is the industry and research community backing that is derived from the use of widely accepted standards based mark-up languages. This open community backing ensures the availability of quality editing tools and repository and reasoning servers along with a large body of published knowledge in the Semantic Web technologies field.

An example of the direct annotation of terms in a legal document is shown in figure 3.

Figure 3 demonstrates the annotation of two of the components of the example citation "AAA v BBB (1998) 23 FAM LR 716, affirmed" given previously in the *Case Report* extract in figure 1. This is achieved using the XML representation of OWL to mark-up "AAA v BBB" as an instance of the *Considered Case* concept that is implicit in the citation and to define "AAA v BBB" as having a *has_considered_value* property of "affirmed" thereby defining "AAA v BBB" a previously considered case that has been affirmed.



```
<!-- An Affirmed Case has a Considered value of affirmed and
     something that has a considered value is a Considered Case
     and
     Affirmed is a type of Considered Value
     An Affirmed Case is a Considered Case with a considered
     value of Affirmed -->

<!-- The individual instance AAAvBBB is
     a Considered Case with a considered value of Affirmed
     IT is therefore an Affirmed Case -->
<ClassAssertion>
    <OWLClass URI="Considered_Case"/>
    <Individual URI="AAAvBBB"/>
</ClassAssertion>
<ClassAssertion>
    <ObjectSomeValuesFrom>
        <ObjectProperty URI="has_considered_value"/>
        <OWLClass URI="Affirmed"/>
    </ObjectSomeValuesFrom>
    <Individual URI="AAAvBBB"/>
</ClassAssertion>
```

Figure 3: Case Report Mark-up.

## 3.3 Semantic Applications

After the conceptual infrastructure of the ontologies and annotated documents are created, the applications layer of the framework can be used to develop applications for purposes including document retrieval, navigation and domain understanding, document management, legal process automation and legal reasoning.

### 3.3.1 Document Retrieval, Navigation and Domain Understanding

Understanding legal concepts and arguments is important in the legal domain and dealing in concepts rather than documents enhances the access and understanding of legal knowledge particularly when the understanding of a legal concept requires information spanning multiple documents. In cases where documents are made available through web pages, HTML provides simple links that have led to fields of knowledge discovery built from the statistical analysis of links into and out of a page, web site or domain. The introduction of conceptual links brings greater complexities to the link relationships and requires a Semantic Web language to define and express the necessary information. Links can be enriched using RDF to hold more properties than the traditional resource locator and it is possible for links to be multifaceted and allow for reasoning over legal concepts by using OWL to add higher level concepts. Potential benefits of this type of annotation are the development of applications that bring about

- a reduction in time searching legal documents

- a reduction in the amount of non relevant material the legal practitioner must process and

- a reduction of the amount of instruction a legal practitioner must give to a computer to retrieve the required information.

There is a view held that future search engines will rely less on ranking of relevance of complete documents and behave more like information management tools (Bontcheva et al., 2006, p143). These engines will analyse and summarize the content of documents rather than listing the documents by relevance and leaving the determination of the documents meaning to the user. Information will be defined by meaning and not statistical coincidence of terms and searching will not be for documents, but for ideas or concepts within documents or across documents. Legal documents such as a *Case Reports* or legislation have structures that deal with multiple concepts of law. The evolution and refinement of legal concepts across judicial processes such as court appeals and reapplication of the law to new facts, reinforces the importance of understanding concepts across documents.

### 3.3.2 Document Management and Automated Processing

A significant problem in the legal publishing domain is changes to documents such as legislation and its contained law, and the changes to legal concepts as case law evolves. This problem is separate to the substantive legal issues dealt with by legal practitioners but is vital to finding and understanding the applicable laws. The repeal of legislation, the affirmation or reversing of court decisions and the jurisdictional boundaries and authority of courts are all examples of metadata contained in legal documents that can be enriched to provide better ways to manage legal documents to assist in better understanding the application of the law. Document management of court cases and legislation is particularly important for the maintenance of the state of the law because laws should not be retrospectively applied by courts and therefore it is essential to know the state of the law at the point in time of the occurrence of some legal event.

Legislation has a number of stages that it must pass before it becomes law. Courts also have a number of steps that a document must pass through in order to comply with court procedures. The annotation of legal documents provides the information required for systems to track the progress of the production of these legal documents.

### 3.3.3 Legal Reasoning

Much effort has been invested by the AI community to develop decision support systems built on different forms of legal reasoning including, case based, rule based and Description Logics (DL) based reasoning. Our approach limits the use of reasoning to the evolution of higher level legal constructs. The OWL language provides for reasoning using Description Logics and more recently, a constrained horn clause language called the Semantic Web Rule Language (SWRL). SWRL is a formal submission to the W3C (Horrocks et al., 2004) for the introduction of horn clause statements that can be used for the expression of rules such as those found in legislative norms. The use of horn clauses to express legal norms is demonstrated by (McCarty, 1989) in his Language for Legal Discourse and also discussed in (Visser and Bench-Capon, 1998) and can be used to express some norms that can not be easily expressed using DL.

The legal concept of *affirmed on Appeal* provides a good example of the difficulty to assert in OWL DL that a case *Appealed* and *Affirmed* in the same hearing is a case *Affirmed on Appeal*. By contrast, the same concept can easily be expressed in SWRL as

$$AffirmedOnAppeal(x) \leftarrow AffirmedIn(x,y) \wedge AppealedIn(x,z) \wedge sameHearingAs(y,z)$$

The reasoning capabilities of OWL DL can be used to classify new concepts and can therefore be used to create hypothetical scenarios that can be constructed using OWL class constructors and classified for querying purposes.

A further benefit of using the OWL language is that the reasoning capabilities it provides eases the annotation process in circumstances where a simple term in the legal document is to be further developed as a more complex legal concept. This benefit is demonstrated using the the previous case citation example shown in the *Case Report* extract in figure 1. The term *Affirmed* is a member of a collection of attributes for determining the impact a *Considered Case* has on a previous court decision along with the terms *Followed*, *Not Followed*, *Applied* and *Reversed*. At face value, these terms are merely a list of discriminating attributes for a *Considered Case* and would be created as instances of an OWL class and be easily marked-up in the legal document. If however the *Affirmed* attribute is to be expanded into a richer concept, it must be created as an OWL class and therefore can not be annotated directly in the document without also creating a specific *AAA v BBB* instance of the new *Affirmed* class.

Figure 4 demonstrates how the OWL language enables the individual instance *AAA v BBB* to be defined as a type of anonymous class having some value
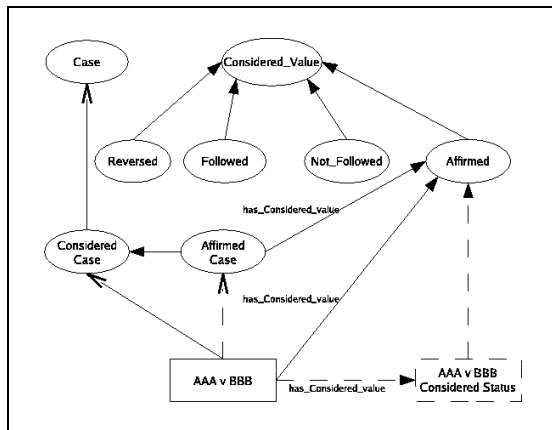
Figure 4: Case Citation Representation.

from *Affirmed_Case* for its *has_considered_value* property.

This enables a reasoner to infer that *AAA v BBB* is an *Affirmed Case* without having to create a particular instance of *Affirmed Case* for *AAA v BBB's has_considered_value* property. The term Affirmed can now be expanded as a legal concept to incorporate information about the jurisdictional hierarchy and relative superiority of the court that heard the case or could include a temporal dimension that indicates the changing authority of a case over time. This example demonstrates the use of new Semantic Web technologies and their associated patterns of best practice (Rector and Welty, 2005) to enrich the concepts found in legal documents in ways previously unavailable.

## 4 CONCLUSIONS

The Semantic Web technologies of RDF and OWL enable existing structured and semi-structured legal documents to be enriched to formally represent the complex and interrelated concepts found within these documents. We have described a framework that employs a bottom-up approach to developing legal ontologies grounded in the terms and structures of the documents and re-using much of the existing mark-up found in todays legal documents as well as the skills of domain experts such as editors, case reporters and legislative drafters. The framework also provides for the direct annotation of these documents by the domain experts and the selective use of the the ontologies in semantic applications. We have also discussed the benefits of this approach and give concrete examples in OWL and RDF. It is considered that such an approach will facilitate more direct and easier legal document annotation, as the ontology concepts are grounded in the existing terms found in the documents themselves.

## REFERENCES

Atkinson, K., Bench-Capon, T., and McBurney, P. (2005). Arguing about cases as practical reasoning. In *ICAIL 2005: Proceedings of the 10th international conference on Artificial intelligence and law*, pages 35–44. ACM Press, New York, NY, USA.

AustLII (2007). Australasian legal information institute. http://www.austlii.edu.au/.

Bentham, J. (1970). *An Introduction to the Principles of Morals and Legislation*. Athlone Press, London.

Biagioli, C., Francesconi, E., Passerini, A., Montemagni, S., and Soria, C. (2005). Automatic semantics extraction in law documents. In *ICAIL 05: Proceedings of the 10th international conference on Artificial intelligence and law*, pages 133–140. ACM Press, New York, NY, USA.

Bontcheva, K., Davies, J., Duke, A., Glover, T., Kings, N., and Thurlow, I. (2006). *Semantic Web Technologies: trends and research in ontology-based systems*, chapter 8, pages 139–169. John Wiley & Sons, Ltd.

Breuker, J. and Hoekstra, R. (2004). Epistemology and ontology in core ontologies: Folaw and lri-core, two core ontologies for law. In *Proceedings of EKAW Workshop on Core ontologies*. CEUR.

Casellas, N., Jakulin, A., Vallb, J.-J., and Casanovas, P. (2006). Acquiring an ontology from the text. In Ali, M. and Dapoigny, R., editors, *Advances in Applied Artificial Intelligence*, pages 1000–1013. Springer.

Hafner, C. D. (1987). Conceptual organization of case law knowledge bases. In *ICAIL '87: Proceedings of the 1st international conference on Artificial intelligence and law*, pages 35–42, New York, NY, USA. ACM Press.

Hart, H. (1961). *The Concept of Law*. Clarendon Press, Oxford.

Horrocks, I., Patel-Schneider, P. F., Boley, H., Tabet, S., and Dean, B. G. M. (2004). SWRL: A Semantic Web Rule Language Combining OWL and RuleML. Technical report, W3C.

Kelsen, H. (1991). *General Theory of Norms*. Clarendon Press, Oxford.

McCarty, L. T. (1989). A Language for Legal Discourse. In *Proceedings of the Second International Conference on Artificial Intelligence and Law*, pages 180–189, Vancouver, Canada.

NYSUC (2007). New York State Unified Court System, law libraries, glossary of legal terms.

Rector, A. and Welty, C. (2005). Simple part-whole relations in OWL Ontologies. Best practices, W3C, http://www.w3.org/2001/sw/BestPractices/OEP/ SimplePartWhole/. Retrieved August 2007.

Valente, A. (1995). *Legal knowledge engineering: A modelling approach*, volume 30 of *Frontiers in artificial intelligence and applications*. IOS Press.

van Noortwijk, K., Visser, J., and Mulder, R. V. D. (2005). Re-usable retrieval concepts for the classification of legal documents. In *Proceedings of the 10th international conference on Artificial intelligence and law*, pages 252–253. ACM Press, Bologna, Italy.

Visser, P. R. S. and Bench-Capon, T. J. M. (1998). A comparison of four ontologies for the design of legal knowledge systems. *Artificial Intelligence and Law*, 6(1):27–57.