# The Open Annotation Collaboration Phase I:
# Towards a Shared, Interoperable Data Model for Scholarly Annotation

Timothy W. Cole, University of Illinois at Urbana-Champaign, University Library and Graduate School of Library & Information Science
Myung-Ja Han, University of Illinois at Urbana-Champaign, University Library

**Abstract**

This paper reports on preliminary outcomes from Phase I of the Open Annotation Collaboration (OAC), discussing them in the context of illustrative scholarly annotation use cases drawn largely from the domain of renaissance emblem studies. The OAC Phase I project sought to address problems of dysfunction caused by too many different, insufficiently interoperable annotation clients and tools through the development and promulgation of a more resource-centric and web-centric standard for making and disseminating scholarly annotations of Web resources. By focusing on representative use cases and an underlying data model of scholarly annotation more consistent with Semantic Web and Linked Data principles rather than on application-specific or interface-specific issues to do with annotation, the OAC seeks to foster annotation sharing and interoperability. Results to date confirm diverse and complex user requirements in regard to the creation and use of scholarly annotations. Nonetheless, a reasonably straightforward and elastic data model is emerging with seemingly good potential to work across a broad spectrum of scholarly annotation use cases and applications. This suggests for Phase II an opportunity for in-depth, domain-specific experiments to further test and refine the initial OAC data model created.

**Introduction**

*Annotating*, the act of associating one piece of information with one (or more) other piece(s) of information, is a core and pervasive practice in the humanities. Annotations are used to create, elaborate, organize and share. Some individual scholars annotate when reading, as an aid to memory, as a way to add commentary, as a way to classify documents. There exists a plethora of annotation clients available to humanities scholars.[1] Some of these tools are general purpose. Others are designed for specific collections or collection types, to address specific user requirements, or to meet a domain-specific need.

As a result (and as has been pointed out by many[2]), scholars have to learn different annotation clients for different content repositories, often have no easy way to integrate annotations made on different systems or created by colleagues using other tools, and often are limited to simplistic and constrained models of annotations suitable for use only in limited contexts. For example,

---

[1] E.g.: Jane Hunter, "Collaborative Semantic Tagging and Annotation Systems," *Annual Review of Information Science and Technology (ARIST)* 43 (2009): 187 – 239; J. Wolfe, "Annotation technologies: A software and research review," *Computer and Composition* 19, no. 4 (2002): 471-97; J. Bradley and P. Vetch, "Supporting annotation as a scholarly tool-experiences from the online Chopin variorum edition," *Literary & Linguistic Computing* 22, no. 2 (2007): 225-41.

[2] E.g.: C. I. Borgman, "Digital libraries and the continuum of scholarly communication," *Journal of Documentation* 56, no. 4 (2000): 412-30; M. Agosti, et. al, "DiLAS: a Digital Library Annotation Service," *Proceedings of Annotation for Collaboration -- A Workshop on Annotation Models, Tools and Practices* (2006), accessed October 8, 2010, http://www.is.informatik.uni-duisburg.de/bib/pdf/ir/Agosti_etal:05.pdf.

many existing tools only support brief unformatted annotation content. Equally problematic, many tools fail to treat annotations as first-class, independent information objects. Such tools instead conflate the storage of the annotation and the target resource being annotated. This approach can frustrate subsequent efforts to directly reference or annotate an annotation in its own right.

This paper reports on preliminary outcomes from Phase I of the Open Annotation Collaboration (OAC). Funded in the spring of 2009 with a generous grant from the Andrew W. Mellon Foundation, this applied research project seeks to address the issues outlined above through the development and promulgation of more resource-centric and web-centric standards for making and disseminating scholarly annotations of Web resources. By focusing first on representative use cases and the underlying data model of scholarly annotation rather than on application-specific or interface-specific issues to do with annotation, the OAC seeks to foster annotation sharing and interoperability. Not surprisingly, preliminary results confirm diverse and complex user requirements in regard to the creation and use of scholarly annotations. Nonetheless, a reasonably straightforward and elastic data model is emerging with seemingly good potential to work across a broad spectrum of scholarly annotation use cases. This suggests for Phase II an opportunity for in-depth, domain-specific experiments to further test and refine the initial data model created.

For larger, higher quality versions of the figures reproduced here, please refer to the *Supplementary Files* section accompanying this article online at http://jdhcs.uchicago.edu

**The Open Annotation Collaboration**

Problems of annotating on the Web cut across disciplines, formats, and research domains. Accordingly a broad range of research perspectives are needed to address these issues. The founding members of the OAC are:

- The Center for History and New Media (George Mason University)
- The eResearch Lab (School of ITEE, the University of Queensland)
- JSTOR
- The Maryland Institute for the Humanities (University of Maryland)
- The Research Library, Los Alamos National Laboratory
- The University Library & Center for Informatics Research in Science and Scholarship (University of Illinois at Urbana-Champaign)

Collectively the members of the Collaboration have committed to three long-term objectives:

- To facilitate the emergence of a Web and resource-centric interoperable annotation environment that allows leveraging annotations across the boundaries of annotation clients, annotation servers, and content collections.
- To demonstrate, through prototype implementations, an interoperable annotation environment in a variety of settings characterized by a range of annotation client/server configurations, content collections, and scholarly use cases.
- To seed widespread adoption of OAC standards in scholarly contexts by deploying robust, production-quality applications conformant with this interoperable annotation environment.

The Collaboration was formed late in 2008 with work on Phase I beginning in the spring of 2009. The Collaboration is currently in Phase II, which began in January 2011.

## Prior Work, Foundational Principles, Motivation

There is of course a large body of past and ongoing work in this domain, but despite a large body of prior art regarding annotation practice, annotation models, and annotation systems, comparatively little attention has been paid to interoperable annotation environments. A few notable efforts in this realm to date:

- RDF-based *Annotea* developed by Kahan and Koivunen;[3]
- Agosti's *Formal Model of Annotations of Digital Content* (Agosti and Ferro 2007);[4]
- Boot's *SANE: Scholarly Annotation Exchange*, based on a model of third-party annotations presented at DH2006 (Boot 2006),[5]
- *OATS: The Open Annotation and Tagging System.*[6]

An analysis of these existing models reveals that on the whole, most have not been designed as Web-centric and resource-centric and/or that they have modeling shortcomings that prevent any existing resource from being the content or target of an annotation and preclude an annotation from being given independent status as a resource itself. To help ensure an advance on prior art, the OAC has articulated a set of guiding principles[7] which inform our work, including:

- The OAC shall focus on interoperability across clients, tools & collections; not on prescribing client interfaces or internal architecture.
- Consistent with prior work, the OAC will model an annotation as a resource linking an annotation body (content) to an annotation target.
- Contrary to some prior work, the OAC will model annotation & annotation body are separable resources with separate identities on the Web.
- Contrary to some prior work, annotation body (content) is not limited exclusively to text types and formats.
- OAC data model must accommodate annotations involving multiple body and/or multiple target resources and specific segments, representations and/or versions of resources in these roles.

---

[3] J. Kahan, M. Koivunen, E. Prud'Hommeaux, and R. Swick, "Annotea: An Open RDF Infrastructure for Shared Web Annotations," *Proceedings of the 10th International conference on the World Wide Web*, Hong Kong, May 2001.

[4] M. Agosti and N. Ferro, "A Formal Model of Annotations of Digital Content," *ACM Transactions on Information Systems* 26, no. 1 (2007), accessed October 8, 2010, http://dx.doi.org/10.1145/1292591.1292594.

[5] P. Boot, "Third-Party Annotations in the Digital Edition Using EDITOR," *Proceedings of Digital Humanities* (Paris: Université Paris - Sorbonne, 2006) 34-35.

[6] S. Bateman, R. Farzan, P. Brusilovsky, and G. McCalla, "OATS: The Open Annotation and Tagging System," *Proceedings of the Third Annual International Scientific Conference of the Learning Object Repository Research Network*, Montreal, November 2006.

[7] Herbert Van de Sompel, Robert Sanderson, Timothy W. Cole, and Jane Hunter, "Open Annotation Collaboration Interoperability Thread: Guiding Principles," accessed October 8, 2010, http://www.openannotation.org/documents/OAC_GuidingPrinciples_20091106.pdf.

- OAC data model defines classes, entities, properties & relationships that facilitate interoperability, but which are also extensible.

With prior work in mind and having articulated principles to guide the effort, the Collaboration turned its attention a range of challenges and questions. These in turn serve as motivation for the effort and provide a standard against which to measure progress:

- Can we describe a broadly useful model of annotation not tied to repository design or type of content being annotated?
- Using this model, can we enable new opportunities for digitally-based scholarship built around annotation & annotation interoperability?
- What are the defining scholarly use cases and can we embed the OAC model in existing applications to demonstrate benefits for these use cases?
- Are there additional benefits to be had by treating annotations as first-class Web Resources?

## The OAC Data Model

A third iteration of the still "alpha" version of the OAC data model to support annotation sharing and interoperability was released in the fall of 2010.[8] This data model provides a method of describing annotations so that they easily can be shared between platforms and across repositories. The model provides sufficient extensibility to support a richness of expression necessary to satisfy scholars' needs while reducing to a simple baseline instantiation to allow for common use cases such as attaching a piece of text to a single Web resource. Consistent with the recommended practices of the Linked Data Initiative,[9] annotations are modeled as a set of connected (HTTP) URI-addressable resources, including one or more *annotation body* resources, i.e. the annotation content or source, and one or more *annotation target* resources. Figure 1 depicts the baseline (simplest) instantiation of the OAC data model represented as a graph of connected resources and properties.

---

[8] "Open Annotation: Alpha3 Data Model Guide," accessed October 11, 2010, http://www.openannotation.org/spec/alpha3/.

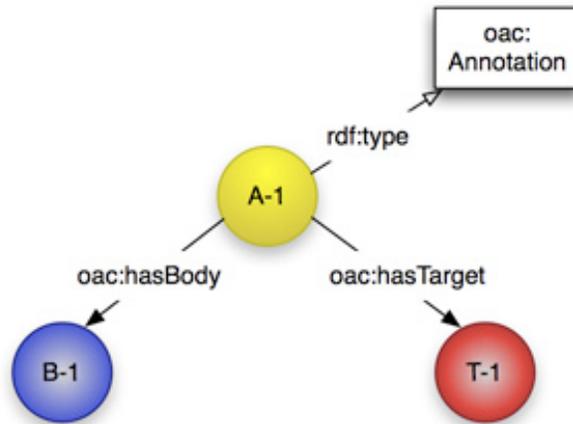[9] "LinkedData," accessed October 11, 2010, http://www.w3.org/wiki/LinkedData.

**Figure 1.** The baseline OAC data model.

While an essential aspect of any annotation is the expression of an *annotates* relationship (at least implicitly) between the body and target, the approach of treating the annotation, the body, and the target as distinct, URI-addressable resources simplifies and decouples an annotation tool implementation from the Web-addressable repository holding content being annotated. Another strength of the OAC data model is that it allows for the annotation body and annotation target to be of any media type. Using this model, the claim can be made that any specific Web resource annotates any other Web resource. Thus a Web-accessible streaming audio (e.g., commentary) can annotate a Web-accessible video clip. Moreover the model allows the annotation, the body, and the target each to be stored separate one from another, and for the annotation, the body, and the target to all have different authorship. This latter feature supports ontological-based annotation, i.e. the annotation of resources using concepts drawn from pre-existing discipline-specific ontology.[10]

Of course the simplest case annotation illustrated in Figure 1 is not sufficient for most scholarly annotation use cases. To leverage existing Semantic Web practices and facilitate interoperability of annotation applications, the OAC data model allows for the expression of additional properties and relationships. These relationships can be attached to any of the annotation, annotation body, or annotation target.

As mentioned above, this allows, at least in theory, separate authorship of annotation and annotation body. More importantly this approach also ensures separate identity for the annotation body (content) and the annotation (the assertion of body annotates target). But while the OAC data model avoids conflating the annotation body with the annotation, the instantiation of the annotation and the annotation body often will occur simultaneously. For example, an annotation can come into existence at the same time as the textual comment which comprises the body of the annotation. In such cases it is convenient to instantiate the annotation body inline, i.e., within the RDF serialization of the annotation description. The OAC data model

---

[10] N. H. Shah et al., "Ontology-Driven Indexing of Public Datasets for Translational Bioinformatics," *BMC Bioinformatics* 10 (2009), Supplement 2 (S1), accessed October 11, 2010, http://dx.doi.org/10.1186/1471-2105-10-S2-S1.

leverages the W3C's *Representing Content in RDF* Working Draft[11] to allow this in a manner consistent with Semantic Web best practice. The *Representing Content in RDF* Working Draft supports the embedding of plain text, XML (including XHTML) and base64-encoded annotation bodies. To allow reference to the annotation distinct from the annotation body, it is still required that the annotation body have a unique URI. To create this unique URI for an inline annotation body, a *urn:uuid* may be used.

Scholarly annotations may also target only parts of resources or may target multiple resources. For instance annotations that compare and contrast resources necessarily involve multiple annotation targets. Scholarly annotations may involve a body or a target that is a specific representation of a resource or a segment or fragment of a particular resource. Scholarly annotations may involve annotation body or target resources that are ephemeral, subject to change, or have other time dependencies. The OAC data model provides features by which these more complex annotations can be described. The use of URIs including fragment identifiers is encouraged. In order to allow for use cases which cannot be described using fragment URIs alone, the OAC data model defines two additional entities, the *ConstrainedTarget* and the *ConstrainedBody*, and two special predicates, *constrains* and *constrainedBy*. These entities and relationships can be used to describe annotations that target (or have as body) a specific segment, representation, or version of a resource.

Figure 2 depicts three graphs illustrating how the model can be extended to describe annotations more sophisticated than the baseline, simplest annotation depicted in Figure 1.
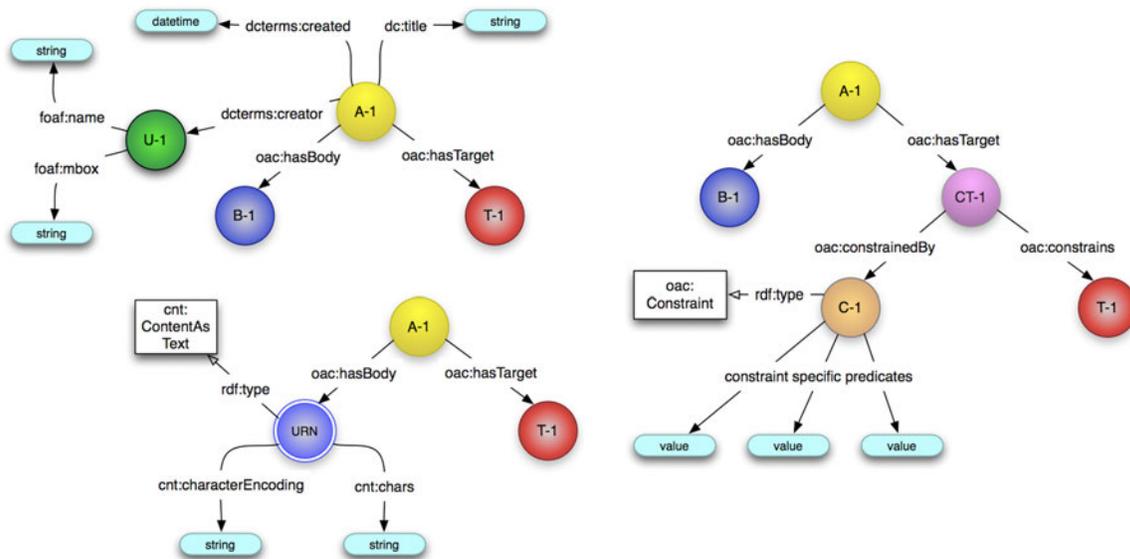


**Figure 2.** More complex use cases represented in the OAC data model.

_____

[11] "Representing Content in RDF 1.0 (W3C Working Draft)," (2011), accessed June 3, 2011, http://www.w3.org/TR/Content-in-RDF10/.

## Use Cases

To inform and support the development of the OAC data model, a range of scholarly annotation use cases were examined. These initial use cases were developed from a review of the literature, augmented by direct discussions with scholars in multiple disciplines. Annotations involving both digital and non-digital resources were examined.

Our initial examination of scholarly annotation use cases has raised other data modeling issues and questions beyond those mentioned above. Many of these questions remain to be resolved. Among them:

- Through HTTP content negotiation it is sometimes possible to dereference a given URI so as to retrieve a preferred representation of a resource, e.g. to retrieve an image in one format rather than another, or a text in French rather than English. The OAC data model provides a means to reference a preferred representation of a given resource as annotation body or target. How often is this functionality required in practice?
- As illustrated above, the OAC data model accommodates annotations involving multiple targets as well as annotations targeting Aggregations described according to the *Open Archives Initiative Object Reuse and Exchange* (OAI-ORE) specification.[12] What factors should implementers consider when deciding which approach to use?
- Work to date suggests that issues of anchor vs. citation may be important in some scholarly domains. A literary scholar may annotate a passage in a novel while viewing a specific digital instance of that novel. While the annotation target can be anchored in the specific digital instance, often the intent is to annotate the passage in multiple digital instances of the novel, possibly spanning editions of the work. Can such intent be expressed using the OAC data model?
- An annotation may both target one resource and reference another resource. How is this use case distinguished from an annotation involving multiple annotation targets?
- Biodiversity annotation use cases examined suggest that there are times when a user may wish to annotate a resource in context, e.g. a scholar may want to say something about an image, but only about the image as it is embedded in a specific Web page. How important a use case is this in other domains, and if important, is it well enough handled by the OAC data model.

### Annotation of Digitized Renaissance Emblems – An Illustrative Use Case

In September of 2009, the University of Illinois at Urbana-Champaign and the Herzog August Bibliothek, Wolfenbüttel received funding from the National Endowment for the Humanities (U.S.) and the Deutsche Forschungsgemeinschaft (Germany) for the *Emblematica Online* project.[13] As part of this project, the University of Illinois is digitizing 424 emblem books from its Rare Book and Manuscript Library collections, including a distinguished subset of about 50 German language emblem books. These books in particular are being analyzed with finer grained elements (e.g., emblem pictura) being made individually addressable and discoverable, while

---

[12] "Open Archives Initiative Object Reuse and Exchange," accessed June 3, 2011, http://www.openarchives.org/ore/.

[13] "Events and Activities, Department of Germanic Languages and Literatures, University of Illinois at Urbana-Champaign," accessed June 3, 2011, http://www.germanic.illinois.edu/news/emblem/.

simultaneously maintaining connection to parent book as context for the more fine grained elements. In combination with several other key Web-accessible emblem collections at HAB and elsewhere, these resources comprise an important corpus for scholarship and in particular scholarly annotation.
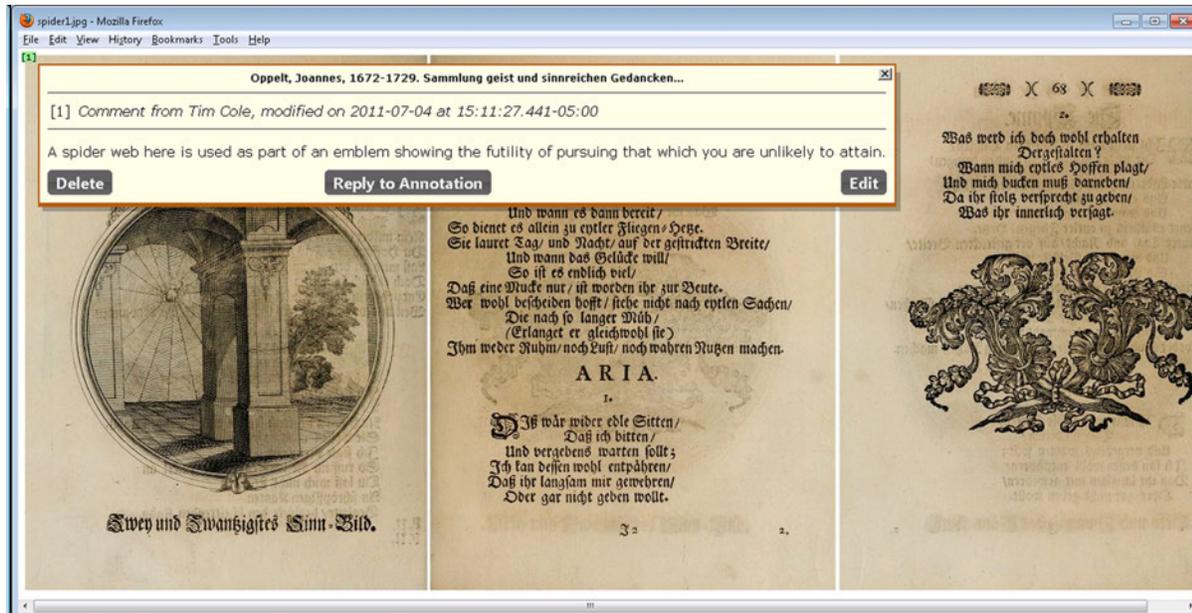


**Figure 3.** Simple annotation of a digitized emblem.

So, for example, it is relatively straightforward imaging a scholar or student studying the use of spider iconography in emblems wanting to annotate an instance of an emblem discovered as shown in Figure 3. This annotation is relatively easy to model, though to do the intent of the scholar full justice it is necessary to constrain the target of the annotation to the region of the image containing the spider web. This annotation is easily handled by the draft OAC data model.

Consider next that this same scholar, or potentially another scholar, having created or encountered annotations of three instances of spider imagery in emblems might then want group these three instances (as depicted in Figure 4) and create an additional, new multi-target annotation commenting on the similarity or differences in the use of iconography in these emblems. Because the OAC data model allows annotations to be identified and addressed separately from annotation targets, this again is handled by the OAC data model. The new multi-target annotation would unambiguously target the three prior annotations, not their annotation bodies (i.e., not the textual comments made about each individual emblem).
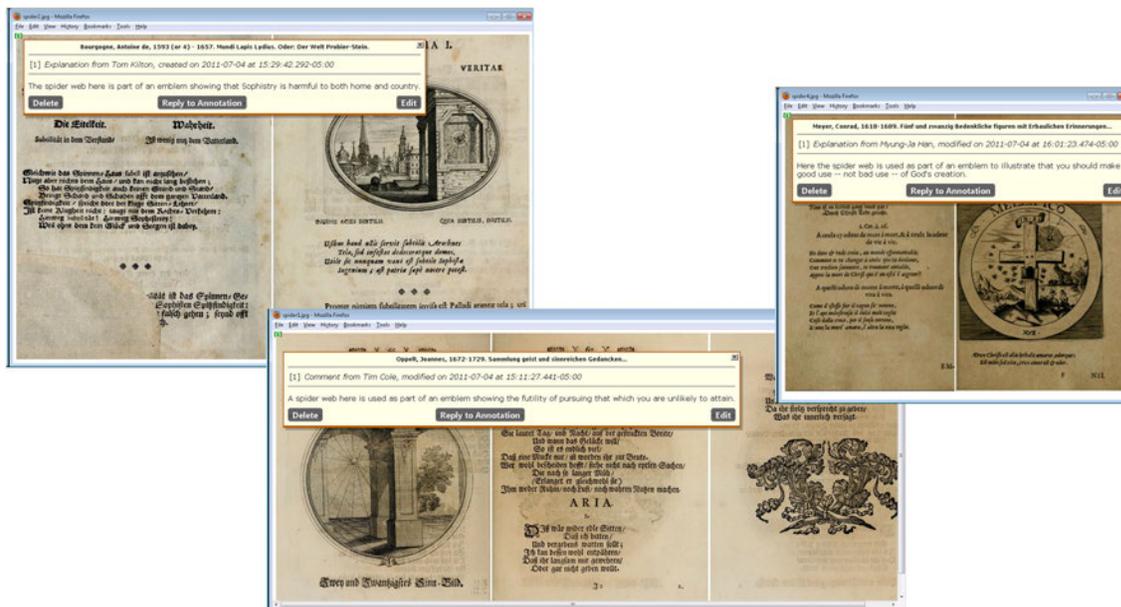
**Figure 4.** An annotation of three previously created annotations.

Consider finally the wholly distinct example depicted in Figure 5. Here we have an emblem pictura to which Iconclass subject headings have been assigned. In essence the assignment of each heading can be treated as an act of annotation, with the body being the heading assigned (each of which has its own identity) and the target being the pictura region of the emblem. For the most part the OAC model does a good job representing these annotations. The OAC separation of annotation from annotation body means it is easy to reference as annotation body the Iconclass heading within its canonical hierarchy. However, the matter of constraint is not entirely straightforward for all of the headings shown. The arms raised and heart symbolism plus fire headings annotations are easily managed by constraining the target to regions of the pictura (though in the former instance, the annotation will have two targets, one for each hand depicted). The proper target constraint for the third heading, 57A8(+4) Gratitude is a bit more nuanced. Here a value is being assigned to the pictura. It can be forcefully argued that while the Iconclass heading is indeed annotating the pictura, it is doing so only in the context of the particular emblem. Pictura were from time to time reused in whole (or more often in part) in different emblems. Sometimes the result was significantly different meanings. It is not yet clear how to properly express, in a machine-understandable way, a constraint that essentially says, "the whole of this image, but only when viewed in the context of this emblem." More research and experimentation is needed.

**Figure 5.** Annotations of an emblem pictura.

## Conclusion

The proposed OAC Data Model has the potential to enable the sharing and discovery of annotations beyond the boundaries of individual solutions and content collections, and hence the potential to allow for the emergence of value-added, cross-environment annotation services. It also has the potential to facilitate the implementation of advanced end-user annotation services that are capable of operating across a broad range of both scholarly and general collections. Furthermore, it has the potential to enable customization of annotation services for specific scholarly communities, without reducing interoperability, and enable more robust machine-to-machine interactions and automated analysis, aggregation and reasoning over distributed annotations and annotated resources. The main elements of the OAC data model are now in place. Through further experimentation during Phase II of the OAC project involving discipline-specific annotation use cases, the model will be refined and (we hope) the major remaining issues resolved. By grounding this work in a thorough understanding of Web-centric interoperability and embedded models implemented by existing digital annotation tools and services, the OAC hopes to create an interoperable annotation environment that will allow scholars and tool-builders to leverage prior tool development work and traditional models of scholarly annotation, while simultaneously enabling the evolution of these models and tools to make the most of the potential offered by the Web environment.

## Acknowledgements

**Bibliography**

Agosti, M. and N. Ferro. "A Formal Model of Annotations of Digital Content." *ACM Transactions on Information Systems* 26, no. 1 (2007). Accessed 8 October 2010. http://dx.doi.org/10.1145/1292591.1292594.

Agosti, M., et. al. "DiLAS: a Digital Library Annotation Service." *Proceedings of Annotation for Collaboration -- A Workshop on Annotation Models, Tools and Practices.* (2006), accessed May 31, 2011. http://www.is.informatik.uni-duisburg.de/bib/pdf/ir/Agosti_etal:05.pdf.

Bateman, S., R. Farzan, P. Brusilovsky, and G. McCalla. "OATS: The Open Annotation and Tagging System." *Proceedings of the Third Annual International Scientific Conference of the Learning Object Repository Research Network* (2006). Accessed October 8, 2010. http://www.cs.usask.ca/~ssb609/files/oats-lornet.pdf.

Boot, P. "Third-Party Annotations in the Digital Edition Using EDITOR." *Proceedings of Digital Humanities* (2006): 34-35.

Borgman, C. I."Digital libraries and the continuum of scholarly communication." *Journal of Documentation* 56, no. 4 (2000): 412-30.

Bradley, J. and P. Vetch. "Supporting annotation as a scholarly tool-experiences from the online Chopin variorum edition." *Literary & Linguistic Computing* 22, no. 2 (2007): 225-41.

"Events and Activities, Department of Germanic Languages and Literatures, University of Illinois at Urbana-Champaign." n.d.. Accessed June 3, 2010. http://www.germanic.illinois.edu/news/emblem/.

Hunter J. "Collaborative Semantic Tagging and Annotation Systems." *Annual Review of Information Science and Technology (ARIST)* 43 (2009): 187-239.

Kahan, J., M. Koivunen, E. Prud'Hommeaux, and R. Swick."Annotea: An Open RDF Infrastructure for Shared Web Annotations." *Proceedings of the 10th International conference on the World Wide Web* (2001). Accessed October 8, 2010. http://www10.org/cdrom/papers/488/index.html.

"LinkedData - W3C Wiki." Accessed October 11, 2010. http://www.w3.org/wiki/LinkedData.

"Open Annotation: Alpha Data Model Guide." Accessed October 11, 2010. http://www.openannotation.org/spec/alpha3/.

"Open Archives Initiative Protocol - Object Exchange and Reuse." n.d. Accessed June 3, 2010. http://www.openarchives.org/ore/.

"Representing Content in RDF 1.0." Accessed June 3, 2010. http://www.w3.org/TR/Content-in-RDF10/.

Sanderson, R. and H. Van de Sompel. "Making web annotations persistent over time." *Proceedings of the 10th Annual Joint Conference on Digital Libraries* (2010). Accessed October 8, 2010. http://dx.doi.org/10.1145/1816123.1816125.

Shah, N. H., C. Jonquet, A. P. Chiang, A. J. Butte, R. Chen, and M. A. Musen. "Ontology-Driven Indexing of Public Datasets for Translational Bioinformatics." *BMC Bioinformatics* 10 (2009) (Supplement 2):S1. Accessed October 8, 2010. http://dx.doi.org/10.1186/1471-2105-10-S2-S1.

Wolfe, J. "Annotation technologies: A software and research review." *Computer and Composition* 19, no. 4 (2002): 471-97.