# Rich Linking in a Digital Library
# of Full-Text Scientific Research Reports

ROBERT B. ALLEN
rba@boballen.info
http://boballen.info

**ABSTRACT**

With interactive full-text documents, there are opportunities to take advantage of the structure in scientific research reports which has not been systematically captured. We develop a novel "model-oriented" approach and suggest how that approach may support the development of a new generation of browsers for research reports from the Public Library of Science (PLoS). With full-text we can implement more targeted linking than is possible with monolithic reports. For instance, data sets can be linked directly to the workflows which describe how they were generated and analyzed. Traditional citations can be enhanced by anchoring them to specific points in a cited text. If a reader follows a citation, the model-oriented structure can be used to generate a summary of the target document related to the topic of the citation. In addition, we propose text pre-processing and file standards to facilitate ingest and use of full-text articles.

## 1 VISION: A MODEL-ORIENTED DIGITAL LIBRARY OF RICHLY LINKED FULL-TEXT ARTICLES

Until recently scholarly digital libraries have been composed of collections of complete research reports. However, this has changed with collections of full-text research reports such as those from PLoS. In these changes, there are opportunities to take advantage of the structure in scientific research reports which previously has not been systematically captured. Based on this structure, rich interlinking can be added within the research articles. This approach could make digital libraries more modular, more like knowledgebase, and more like composite [16] and adaptive hypertexts. Wikipedia occupies a point in this design space and some of its features may be common in modular digital libraries in the future. We expect to go well beyond the level of structured interaction provided by Wikipedia.

## 2 TOWARDS A FULL TEXT DIGITAL LIBRARY AND ASSOCIATED END-USER WIDGETS

### 2.1 Implementation Details

PLoS is a full-text Web-based publisher of seven journals in the areas of biology and medicine with several thousand articles online. PLoS allows unrestricted use of the use of its publications with attribution. For this study, a PLoS article by Zhai et al. [26] was picked essentially at random. The XML version of the article was downloaded and processed with a Java program which extracted and organized the XML. The XML formatting was suboptimal for efficiently generating fine-grained links. When a set of citations with sequential numbers was specified, only the range of citations was provided in the XML files. Thus, the missing citation numbers were interpolated. Some features such as the cross referencing of figures was rather complex. Moreover, in checking other articles, especially those from other PLoS journals, some differences in the XMLSchema formats were noted.

The conceptual structure consisting of entities and flows were specified with Java classes. These were coded by hand and mapped to the text. Entities had attribute dimensions which included its state and potentially state changes. The Java class clone function could be used because several versions of very complex entity instances such as *drosophila* with small changes in their genetics were included as part of the experimental manipulations reported in the article. This function needs further optimization and should be standardized. The model-based version of the research report parallels the text and would support flexible navigation. Flows implemented the "causal" interactions between the entities. There were method flows which were triggered by the researcher and conceptual model flows which described

the process being studied. For a more complete description see Section 3 and [7]. Eventually, we should save the collection of preprocessed files so they don't have to be recreated each time an article is loaded.

## 2.2 Browser for Personalized Presentation of Full-Text Articles

A text browser applet was implemented as an extension to the text management program. User interaction widgets [6] were added beyond those available for the HTML version of the articles such as a multi-level table of contents (cf., [12]). More interestingly, the browser gives the users a choice of styles for displaying the reference-list. The reference list can be toggled between presentation in order of appearance (PLoS standard) and alphabetical order (Figure 1).
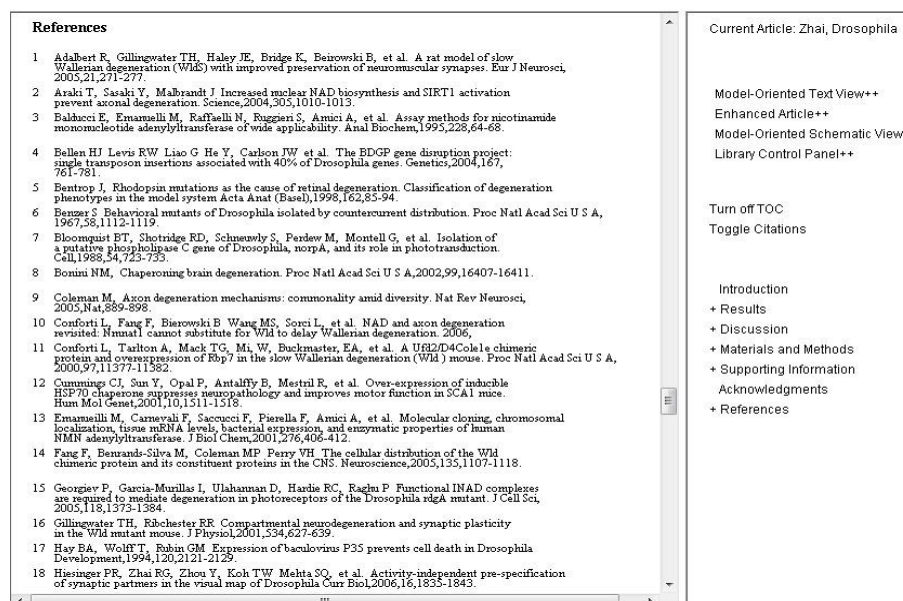


**Figure 1:** Interactive Reference List Widget. The control button toggles the presentation between ordering the references in order of appearance and alphabetically by the first author's surname.

## 3   PROCESS MODELS FOR SCIENCE

The widgets described in the previous section are useful but take only partial advantage of the structure in research reports. Here, we adopt low-level process-based models related to those used in software engineering and business process engineering.

### 3.1 Building on Low-Level Process Units

[8] developed a state-based approach to modeling causation. Its building blocks are entities which have attributes and attributes which have states. There is also a causal relationship when the state of one entity triggers a state change in another entity. In addition, there are variations and extensions of the basic model. For instance, two entities may interact to form new entities as would be the case for the formation of water molecules. The low-level elements can be assembled into extended chains and meshes which incorporate several processes and the models can be expanded (or collapsed) to show (or hide) details.

Workflows are explicit chains which act as a unit and are usually triggered as a unit. As part of the procedure in scientific research, they are usually triggered by a researcher and are used for data collection and analysis. Capturing these should be useful for preservation [1, 2, 21].

### 3.2 Process Models of Scientific Phenomena

It is widely agreed that models are integral to science (e.g., [18, 24]) with great variances in the specifics. Some models are structural (e.g., the Bohr atom) and others describe processes and pathways. There are implicit processes associated with the structural models. Our models can be thought of as describing causal processes (e.g., [25]).

Scientific research may address several different components of the models. It may seek to identify the characteristics of the entities which fit the models including their attributes and states. In some cases, this means finding physical properties which match the models and in other cases, it may require making adjustments to the models. Alternately, research may try to find new relationships among entities. This may involve defining the details of the properties or perhaps in finding abstractions or generalizations of the processes.

The low-level processing elements emphasize qualitative processes. For example, it is natural to say that hydrogen and oxygen molecules interact to form water without including a full quantitative description of how that happens. Likewise, we may say that smoking causes cancer but that is short hand for the full details of that process. While quantitative relationships and feedback loops can be important, we focus on qualitative relationships. This emphasis is consistent with work from psychology which suggests that much human reasoning is qualitative (e.g., [14]). It is also consistent with software engineering and business process engineering. Although we focus on qualitative processes, the approach can be extended to cover quantitative processes. We allow for systems which have complex interactions among the components by treating the system as an entity that can be decomposed into underlying low-level processes (see [4]).

## 4 APPLYING PROCESS MODELS TO RESEARCH REPORTS

There are many opportunities to introduce semantics in scientific publishing (e.g., [11, 22]). Our approach is to separate in the research report the model of the processes under investigation from other aspects of the research such as the experimental procedures and analyses. A structured framework for specifying those models is better than textual descriptions in terms of anchoring and indexing descriptions of the research. Moreover, a flexible approach to modeling based on simple processing units (Section 3.1) will generalize across domains.

These simple process units can be linked to form more complex models and workflows. They may allow more focused browsing of the key points or provide links to contextual material for readers who are unfamiliar with the nuances of a field. Our approach may be distinguished from the studies of scientific research reports which are based primarily on modeling discourse and argumentation (e.g., [10]). Those approaches generally model claims and the evidence used to support those claims. By comparison, our approach focuses on processes rather than claims.

Zhai et al. [26] provided a test case for our approach. Its Introduction is consistent with Swales CARS model (see Section 4.4). The motivation for the research is explained and potentially relevant entities and processes are discussed. From those components and from other evidence, the researcher develops conceptual models. Experiments are then conducted to contrast or extend those models until the possible alternatives are ruled out. Essentially, the conceptual models function as working hypotheses although they depend on the quality of the literature review. The models may evolve as the research is conducted, and the researcher may not have a strong preference among them. Our detailed modeling of these processes may shed light on role of hypotheses and deduction in how science is practiced.

### 4.1 Systematic Descriptions of Research Designs and Methods

A research design is a composite of high-level abstract workflows (sequences of manipulations and measurements) and of tests (Section 4.2) for separate conditions. It might be implemented as a parameterized extension of the notation developed by [9].

A research method is a more specific composite workflow which incorporates the procedure to select test subjects, the assignment of those subjects to conditions, the research design, the conceptual model(s), the manipulation workflow, the measurement workflow, and the data analysis. The term "manipulation" sometimes refers to all of the actions of the researcher on experimental subjects and sometimes only to the differential actions which distinguish between the conditions. Each of these workflows can be specified with the details of some dependent on others'. Potentially, a compiler could validate them.

### 4.2 Measurements, Data Sets, and Analyses

It is common to suggest that data files could be coordinated with a publication's text (e.g., [2, 20]). In our approach, there is a particularly close link between the data and model specification. The workflow models would include the timing and procedures for data collection. Beyond that, the research design provides the logic in terms of the goals of the research. Thus, tests such as the following can be defined [6], where NAD is one of the proteins manipulated in the research:

if ((NAD = =LOW) && (Degeneration = = LOW)) then {prefer ConceptualModel1;}
else if ((NAD = =LOW) && (Degeneration = = NORMAL)) then {prefer ConceptualModel2;}

The research method is instantiated for each test subject. Traditional approaches may save only averages for an entire condition. With the model-oriented approach, it is possible to save several data points for each subject linked directly to the research method instantiation such as the experimental condition, environmental factors (time, research name, etc.), check on manipulations, the instruments setting, and the primary data points. It would also be possible to incorporate several features that are not normally included with the data sets such as links to electronic laboratory notebooks and links to the description of the instruments used, the instruments setting, and instruments calibration.

### 4.3 Summaries Based on the Model-Oriented Descriptions

Structure can be useful for guiding summaries. In our model-oriented approach, the structural descriptions might define which tests were critical for the conclusions reached by the paper and those which had no direct impact on the results. Presumably a summary would focus on the former rather than the latter.

One scenario in which summaries might be useful would be to provide an overview to a reader who was following a citation from another paper. Traditional citations link from a specific point in one article to a second article without a specific anchor in the second article; with full text and model-oriented indexing it should be possible to anchor a citation to a specific point in the cited paper. In some cases, there may not be a specific landing point but multiple points or several points linked by a narrative. Indeed, a summary could be generated about the cited paper from the perspective of the landing point for a reader who followed the citation link. A recent report by Zhai [27] extended the paper we examined above [26] and cited several contributions in the earlier work. With our approach, explicit links could be made to the exact point in the original paper where those details were introduced. For instance, there is a citation in [26] to [27] which refers to neurodegeneration induced by intense light. This effect is not mentioned at all in the abstract of the original paper and appears only indirectly in its table of contents. With a specific anchor in [26], a reader of [27] could easily find details about that point. With a model-oriented structure a context-specific summary could be generated. Figure 2 illustrates a set of options in the browser for presenting such a local summary.
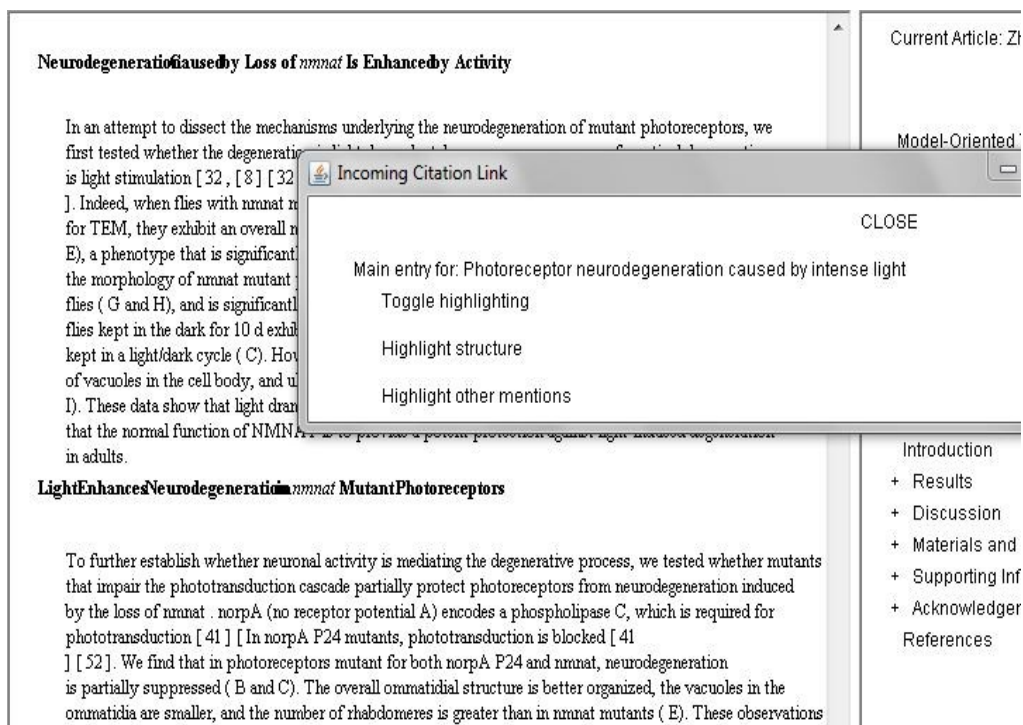
**Figure 2:** A pop-up frame can provide context for a reader of later publications about a specific anchor point. In this case, the links have been generated manually but potentially they could be derived from model-oriented constructs.

### 4.4 Higher-Level Modular Structure

Swales [23] investigated the genres of scientific research reports and emphasized IMRD (Introduction, Methods, Results, Discussion) as the most common structure. Other researchers have described variations of IMRD. Harmsze et al. [17] provide finer-grained components with flows between them.

Swales examined the Introduction in particular depth and described its role as Creating a Research Space (CARS). Swales spent less time on describing the nature of Methods and Results. Our modeling approach may provide tools for examining Methods and Results. Zhai et al. [26] has a particularly complex high-level structure for Methods and Results. In the first half of the "Results" section, it describes the development of several novel methodological techniques for preparing a sample of drosophila with particular characteristics. Then, in the second half of the Results section, it applies the techniques to investigate a set of research questions. We have termed these halves activity blocks and research blocks respectively but they are similar in executing a workflow and reporting the outcome of the manipulation. The main difference between the development of the methods and the research question investigations is whether the investigator is confident about the outcome.

As a further check on the generality of model-oriented approach, we selected a second paper from PLoS essentially at random. This was Gosby et al. [15] from *PLOS ONE*. It had a strikingly different high-level structure from Zhai et al. There was one complex research setting with a between-subjects design in which a large number of measurements were made. Those measurements were then analyzed to explore a number of different conceptual models. Thus, it was challenging to identify distinct blocks for this study.

To handle extended sets of experiments such as found in Zhai et al. [26], we might extend the highly interwoven model of Harmsze et al. [17] with a dynamic flow-control mechanism. To handle studies like Gosby et al. [15] where several questions are investigated in a complex scenario, we might incorporate multiple threads as part of a single block.

### 4.5    Toward a Model-Oriented Digital Library of Scientific Research Reports

A collection of full-text research reports even without full model specifications could be enhanced over the PDF or HTML versions current available. For instance, a browser interface could be deployed, two-way citation links could be implemented, and supplemental pages added. With rich modeling, the model components would need to be kept consistent across articles. Eventually, the model-oriented form may be considered primary to the text version. We plan to save the collection of preprocessed files along with fine-grained models of the content as described below.

## 5    MODEL ORIENTED HISTORY LIBRARIES

Beyond science, there are vast collections of historical materials online. Because of the large amounts of material, organizing them to support access is crucial. Similar to model-oriented scientific research reports, we are developing a model-based fabric of causally related historical events. This builds on a series of studies in which we linked events in timelines to enable them to be more interactive. For instance, in [3] we developed a timeline-like interface to allow exploration of different explanations for the causes of the American Civil War.

In our current work on history, we employ frame semantics [13] to describe events. Specifically, we are using frames from the FrameNet research project as a type of controlled vocabulary. This enables us to capture text from explanatory and narrative histories. While FrameNet frames have been used extensively for marking up linguistic corpora we do not believe they have been applied to information organization. Verb frames are generally consistent with the low-level "causal" state-change models described above for science. The slots of the FrameNet frames provide attribute dimensions for the state changes. Frames can also be applied to modeling scientific research which could be seen as defining and linking new conceptual frames.

History is messier than science. Claims about an event or about the relationship of several events seem best modeled by argumentation systems. Thus, our approach for a model-based view of history adds discourse elements to the event model. Potentially, our modeling will improve access to full-text historical materials such as collections of OCR'd historical newspapers. It might also support services such as linking of museum artifacts to events, developing model-oriented biographies, and linking footnotes in historical analyses with some of the same techniques we have described above for linking citations to scientific research reports.

## 6    CONCLUSION

We have described a browser for interacting with full-text research reports and developed a framework for model-oriented research reports based on low-level processes. Beyond individual research reports, the approach can be extended to develop model-oriented digital libraries. In addition to PLoS, several other publishers currently publish full text of their research reports and several additional publishers should be able to provide full text if requested. Furthermore, our approach of using low-level processes to support model-oriented research reports may be coordinated with libraries of biological process descriptions from the Systems Biology Markup Language (SBML) [19].

Most previous work in library science dealing with information organization has focused on entire documents. There is a need for a new approach to information organization, to develop standards for the description and linking of full text. The goal would be to support interactive exploration, beyond

indexing to support finding information in a rich full-text collection.  In this paper, we have explored using process models and to provide that new approach.

## 7    REFERENCES

1.  Allen, R.B., Using Information Visualization to Support Access of Archival Records. *Journal of Archival Organization*, 3(1), 2005, 37-49.
2.  Allen, R.B., Highly Structured Scientific Publications. *ACM/IEEE Joint Conference on Digital Libraries*, 2007, 472. doi:10.1145/1255175.1255271
3.  Allen, R.B., Visualization, Causation, and History, *iConference,* 2011, doi: 10.1145/1940761.1940835
4.  Allen, R.B., Model-Oriented Scientific Research Reports, *D-Lib Magazine*, May 2011. doi:10.1045/may2011-allen
5.  Allen, R.B., Coordinating Concepts and Discourse in Model-Oriented Research Reports, *ICADL*, Oct. 2011, 392-394, doi:10.1007/978-3-642-24826-9_54
6.  Allen, R.B., Weaving Content with Coordination Widgets, *D-Lib Magazine*, Nov 2011 doi:10.1045/november2011-allen
7.  Allen, R.B., Supporting Structured Browsing for Full-Text Scientific Research Reports, arXiv, 2012, http://arxiv.org/abs/1209.0036
8.  Allen, R.B., Wu, Y.J., and Jun, L., Interactive Causal Schematics for Qualitative Scientific Explanations, *ICADL*, 2005, 411 -415. doi: 10.1007/11599517_50
9.  Campbell, D.T. and Stanley, J.C., *Experimental and Quasi-Experimental Designs*. Rand-McNally, Chicago, 1966.
10. de Waard, A., Breure, L., Kircz, J.G., et al., Modeling Rhetoric in Scientific Publications, 352-356. In *Current Research in Information Science and Techn*ology, 2006.
11. de Waard, A. and Kircz, J.G., Modeling Scientific Research Articles - Shifting Perspectives and Persistent Issues, *ELPUB2008. Conference on Electronic Publishing*, 2008, 234-245. http://thinkubator.ccsp.sfu.ca/wikis/PUB800/Home/uploads/dewaard-elsevier.pdf
12. Egan, D.E., Remde, J.R., Gomez, L.M., Landauer, T.K., Eberhardt, J., and Lochbaum, C.C., Formative Design Evaluation of SuperBook. *ACM Transactions on Information Systems* (*ACM TOIS*), 7, 1989, 30-57. doi: 10.1145/64789.64790
13. Fillmore, C., Frame Semantics and the Nature of Language, *Annals of the New York Academy of Sciences: Conference on the Origin and Development of Language and Speech*, 1976, 280, 20-32. http://www.icsi.berkeley.edu/pubs/ai/framesemantics76.pdf
14. Forbus, K.D., Qualitative reasoning, *CRC Handbook of Computer Science*, 1996. CRC Press, Boca Raton, FL, 715-733.
15. Gosby, A., Conigrave, A.D., Lau, N.S., Iglesias, M.A., Hall, R.M., Jebb, S.A., Brand-Miller, J., Caterson, I.D., Raubenheimer, D., and Simpson, S.J., Testing Protein Leverage in Lean Humans: A Randomised Controlled Experimental Study, *PLOS ONE*, 2011. doi:10.1371/journal.pone.0025929.g001
16. Halasz, F. G., Reflections on "Seven Issues": Hypertext in the Era o the Web. *ACM Journal of Computer Documentation* 25(3), 2001, 109-114. doi: 10.1145/507317.507328
17. Harmsze, F.A.P., van der Tol, M.C., and Kircz, J.G., A Modular Structure for Electronic Scientific Articles, 2008, http://www.science.uva.nl/projects/commphys/ papers/infwet/infwet.html
18. Hestenes, D.  Notes for a Modeling Theory of Science, Cognition, and Instruction, In E. van den Berg, A. Ellermeijer & O. Slooten (eds.) *Modelling in Physics and Physics Education*, (U. Amsterdam 2008), http://modeling.la.asu.edu/r&e/Notes_on_Modeling_Theory.pdf
19. Hucka, M., Bergmann, F., Hoops, S., Keating, S., Sahle, S., Schaff, J., Smith, L., and Wilkinson, D., *The Systems Biology Markup Language (SBML): Language Specification for Level 3 Version 1 Core.* http://www.sbml.org/

20. Hunter, J., Scientific Publication Packages – A Selective Approach to the Communication and Archival of Scientific Output, *International Journal of Digital Curation*, 2008, 1(1), 33-52, doi:10.2218/ijdc.v1i1.4
21. Mayer, R., Rauber, A., Neumann, M.A., Thomson, J., and Antunes, G., Preserving Scientific Processes from Design to Publications, *Theory and Practice of Digital Libraries*, 2012, 113-124, doi: 10.1007/978-3-642-33290-6_13
22. Shotton, D., Semantic Publishing, *Learned Publishing*, 22, 2009, 85-94. doi:10.1087/2009202
23. Swales, J.M., *Genre Analysis: English in Academic and Research Settings*, Cambridge University Press, Cambridge UK, 1990.
24. Thagard, P., *Conceptual Revolutions*, Princeton University Press, Princeton, 1992.
25. White, B., ThinkerTools: Causal Models, Conceptual Change, and Science Education. *Cognition and Instruction*, 10(1), 1993, 1-100. http://www.jstor.org/stable/ 3233779?origin=JSTOR-pdf
26. Zhai, G., Cao, Y., Hiesinger, P.R., Mehta, S.Q., Schulze, K.L., Verstreken, P., Zhou, Y., and Bellen, H.L., Drosophila NMNAT Maintains Neural Integrity Independent of its NAD Synthesis Activity. PLoS Biology, 4(12). 2006, doi:10.1371/journal.pbio.0040416
27. Zhai, G., Zhang, F., Hiesinger, P.R., Cao, Y., Haueter, C.M., and Bellen, H.J., NAD Synthase NMNAT as a Chaperone to Protect Against Neurodegenration. *Nature*, 452(7189), 2008, 887-891 (Letter). doi:10.1038/nature06721