



## ConceptsDB: A New Database Engine and Dataset Model

### Overview

Linguistic Technology Systems (LTS) is developing a new database engine called **ConceptsDB**, which is based on hypergraph data modeling. This new database engine features a multi-paradigm data representation strategy, informed by contemporary research into **AI** and information semantics and permitting theoretical frameworks such as Conceptual Spaces, Conceptual Role Semantics, and Petri Nets to be applied toward the computational modeling of information spaces. In addition, **ConceptsDB** prioritizes application development — particularly native/desktop-style applications. Data about application/session/user state can be directly stored in the database, so that **ConceptsDB** can provide a convenient scaffolding for implementing desktop software. Hypergraph-based data modeling ensures that information needed by the application can be readily marshaled between run-time formats and whichever persistent/serial/**GUI** representations are necessary for application-level capabilities.

LTS is currently working on a prototype of **ConceptsDB** which is part of our “**ConceptsDB** Application Framework” (**CsAF**). This framework is oriented to scientific computing and scientific software. To be specific, **CsAF** database instances can be employed at several sites within a scientific-computing and/or data-sharing platform: as an embedded component of published data sets; as a cloud service managing data and/or publication repositories; or as a tool for individual or institutional users to track data sets and publications. Our “**MOSAIC**” code libraries (“Multi-Paradigm Ontologies for Scientific and Technical Publications”) encapsulate functionality applicable to using **ConceptsDB** in publishing contexts. A **MOSAIC** Portal is a data/publication repository where users can query and download resources following a protocol associated with **ConceptsDB**. The “**MOSAIC** Dataset Explorer” (**MdsX**) libraries are designed for applications which interoperate with data/publication archives (including but not limited to **MOSAIC** Portals). Individual data sets may also embed a version of **ConceptsDB** — specifically, “**DigammaDB**” (**QDB**), which is a light-weight, self-contained **ConceptsDB**-compatible engine suitable to be distributed as source code within a data set/code repository.

To help researchers deploy data sets with customized desktop software, LTS has developed a “Dataset Creator” (**dsC**) designed to be used with **QDB**. This database engine and data-modeling technology (**MdsX**, **QDB**, and **dsC**) form a trio of libraries, which are helpful for the implementation of published data sets via special-purpose applications that are customized for viewing and examining information specific to the research methods and paradigms from which a data set originates; and for the implementation of tools to acquire and keep track of data sets and their associated publications.

Collectively, this trio of libraries form the basis of a **CsAF** implementation — i.e., an Application Framework centered on **ConceptsDB** — which prioritizes deploying and accessing scientific data sets. The information and assets included within a data set, and modeled via **CsAF** components, can span several different requirements, including:

1. Data sets themselves — these may include raw data files and/or code for accessing this data as well as demonstrations/implementations of algorithms for analyzing or otherwise processing the data (or data having a similar format/structure or scientific background);
2. Full-text publications — potentially including both machine-readable and human-readable (e.g. **PDF**) formats, annotated so as to link parts of a text document with corresponding parts/elements within its associated data set(s);
3. Systematic descriptions of research methods, protocols, and (if applicable) equipment, such as may be formally expressed by standards like **MIBBI** (Minimum Information for Biological and Biomedical Investigations), **BIOCORDER** (see Footnote 3), or **PANDORE** (a bioimaging library that contains an “Image Processing Objectives” Ontology);

4. Applications for viewing data sets — applications which may be in the form of pre-existing software, or custom-built for an individual data set.

The **MOSAIC** Dataset Explorer (**MdsX**) libraries include code for interoperating with data structures documenting each of these aspects of scientific resources, which may involve querying data sets themselves or querying **APIs** of scientific corpora where data sets and publications are hosted. **MdsX** is built around what we call a "Scientific Data Repository Model" (**SDRM**) which is concretized separately for each corpus or repository connected with a given **MdsX** application. **SDRM** is divided into distinct concretizations, or "modules," organized around the **APIs** and data profiles of individual scientific portals/repositories. For this reason, LTS seeks to collaborate with organizations who maintain scientific portals so that we can make open-access **SDRM** modules available to the general public, targeting those specific scientific resources.

The following sections will describe **MdsX**, **dsC**, and **QDB** in greater detail as well as provide more information about **SDRM** modules.

## Part I: MOSAIC Components

In recent years — as publishing has become more "digitized" — many online platforms have emerged for scientists to share and publish their research work, including both raw data and scientific papers. These portals generally provide an index/search feature where readers can locate research based on the name of the author, title of the publication, keywords, or subject matter, along with a document viewer where readers can view the abstract of each publication (and sometimes full text of the article) and other publication details (such as co-authors, date of publication, bibliographic references, etc.). However, these portals offer only limited features for interactively exploring publications, particularly when it comes to examining data sets and/or browsing multimedia assets associated with publications, such as audio, video, interactive diagrams, or **3D** graphics.

The **MOSAIC** system is envisioned, in its totality, as a suite of code libraries developed by LTS allowing institutions to host publication repositories — and for individual users to access publication repositories — free of the limitations of existing scientific document portals. In the immediate future, LTS is focused on the more modest goal of implementing software to access existing corpora (via **MdsX**) and to create individual data sets (via **dsC**). With **MOSAIC**, each publication can be linked to a supplemental archive that contains information about the author's research methods, data sets, and focal topics. If desired, these archives could include machine-readable representations of full publication text, to support advanced text-mining techniques across the repository. The supplemental archive can be explicitly linked to publications within their host repository, or it may be maintained in a decentralized manner external to the hosting platform.

Using **MOSAIC**, developers can implement a hosting platform and/or a client platform for this supplemental material, in addition to the publications themselves, where the platform provides software enabling readers to browse and access supplemental archives and their data sets and/or methodological descriptions.<sup>1</sup> **MOSAIC** can be customized for different subject areas by incorporating domain-specific ontologies into document-search features as well as machine-readable document-text annotations.

**MOSAIC** is designed so that the software to access publications and publications' supplemental archives can be embedded in scientific-computing applications via **MOSAIC plugins**. This ensures that publications and data sets can be examined interactively with the same software that scien-

---

<sup>1</sup>A **MOSAIC** publication repository or "portal" is a structured collection of data and documents which can be hosted via web servers (including fully encapsulated and containerizable cloud services) using technology related to **ConceptsDB**. **MOSAIC** repositories can serve as general-purpose portals, hosting academic papers covering a broad range of topics. Alternatively, **MOSAIC** repositories can serve as targeted portals hosting papers focused on more narrowly focused subject domains. Organizations can utilize **MOSAIC** internally as the basis of a private document management system (**DMS**), or to provide (public) open access to complete publications, including human-readable and machine-readable full text and all supplemental content, with no restrictions. Excepting sensitive/private information which might be provided via a dedicated component within the portal with specific features for authors, editors, and administrators. Any project which *does* restrict access to some part of any document requires a commercial **MOSAIC** license.



tists employ to conduct research. **MOSAIC** also introduces a *structured reporting system* (**MOSAIC-SR**) for describing research/experimental/lab methods/protocols. Supplemental archives, plugins, and the structured reporting system are outlined below.

## Supplemental Archives

**MOSAIC** supplemental archives are additional resources paired with **MOSAIC** publications. In general, supplemental archives may include raw data, descriptions of research methods, or annotated data linked to publication texts. Each supplemental archive could have any of the following resources:

1. Interactive versions of the publication, with annotations indicating important concepts and phrases, perhaps aggregated into a "glossary" defining technical terms;
2. Machine-readable representations of document texts, with special-purpose character encodings designed to facilitate Text and Data Mining (**TDM**);
3. Structured files containing raw data discussed in the publication, along with interactive software allowing scientists to access and reuse the data;
4. Detailed reports of the author's research methods and experimental setup and/or protocols, conforming to the relevant standards with respect to the publication's subject classification — for instance, the "Minimum Information for Biological and Biomedical Investigations" (**MIBBI**) includes about 40 specific standards for different branches of biology and medicine;
5. Representations of analytic methods and algorithms underlying the research findings, which are provided directly via computer code or indirectly via formal descriptions of computational workflows;
6. Self-contained computer software which demonstrates code that the author developed and/or used for analyzing/curating research data; and
7. Multi-media assets such as audio or video files, annotated images, **3D** graphics, interactive **2D/3D** plots and diagrams, or other kinds of non-textual content which needs to be viewed with special multi-media software.

The contents of a supplemental archive will be different for different publications, depending on whether the archive contains specific raw data or just a summary of the methods used to obtain the raw data. In the former case, a typical archive will include a *Dataset Application*, or a semi-autonomous software component allowing researchers to study and visualize the data set, typically provided in turn via data files that are also located in the archive. The Dataset Application will, in general, provide both a visual interface for raw data and code libraries for computationally manipulating this data; typically raw data files will encode serialized data structures that can be deserialized by functionality (e.g., **C++** classes) included in the application code (Dataset Applications are discussed further in Part II, Page 7). The **MOSAIC** Dataset Explorer (**MdsX**) libraries can be adapted to whatever specific kinds of information must be serialized for a given publication, providing an implementational starting-point for Dataset Applications.

## MOSAIC Plugins

As an alternative to fully self-contained and autonomous Dataset Applications, **MdsX** explorers can also be deployed as plugins for existing scientific software, allowing publications and supplemental archives to be read and examined within computer software associated with each publication's subject matter. For example, articles about chemistry could be read within **IQMOL**, a molecular visualization program; articles about cellular biology and bioimaging could be read within **CAPTK** (the Cancer Imaging and Phenomics Toolkit), an image-processing application; articles discussing novel computer code might be read within **QT** Creator, an Integrated Development Environment (**IDE**) for programming; etc. The advantage of accessing a publication and a supplemental archive within actual scientific software is that it allows research work to be understood, evaluated, and reused within the computing environments which scientists typically use to conduct professional research. This is different than existing science publication portals, which generally rely on web



browsers to access such supplemental materials as data sets and multimedia files — in this standard setup the software ecosystem wherein readers examine published research is fundamentally separate from the software platforms where research is actually conducted. This example demonstrates how existing portals are limited in their ability to share research in rigorously reusable and replicable ways, motivating MOSAIC's contrarian design.

Embedding presentations of their research within existing scientific software has the added benefit for authors of making their work more practically and immediately useful for the scientific community. Such presentations establish the computational framework for pragmatically deploying their techniques in real-world scientific contexts, accelerating the pace at which research work can be translated to concrete scientific (and clinical/lab/experimental) practice. As one example, research involving novel image analysis techniques could be packaged so as to target a MOSAIC plugin for bioimaging software such as **CAPTK**, so that readers could actually run the author's code as a **CAPTK** module.<sup>2</sup> This is important because such functional assessment and adoption of novel contributions is harder to carry out if a body of research work is described indirectly within article text, as opposed to being concretely implemented within a specific scientific application.

Another benefit of using plugins to access supplemental archives is that the host application will usually provide more sophisticated multimedia and data visualization capabilities, compared to static **PDF** images or even interactive web portals. Publishers have begun to develop online platforms for browsing research papers in conjunction with multimedia content such as interactive diagrams and **3D** graphics — a physical model of a protein or a chemical compound, for instance, can be viewed online via **WEBGL**; such graphics could even be embedded as an "**iframe**" within an **HTML** version of the publications. Publishers consider this to be cutting-edge technology. However, the same molecular **3D** model, when viewed in **IQMOL**, can be enhanced with many additional visual features, representing bonds, orbitals, torsion angles, etc. The multimedia experience of exploring chemistry data in custom software like **IQMOL** is therefore much greater than the experience of generic web multimedia, which means that the scientific software is a better forum for showcasing novel research.

## MOSAIC Structured Reporting (MOSAIC-SR)

The MOSAIC structured reporting framework (MOSAIC-SR) includes tools to help authors develop interactive presentations supplementing academic documents, and specifically to use supplemental archives to document how their research has been conducted. With MOSAIC-SR, authors can implement or reuse code libraries that report on research/experiment methods, workflows, and protocols. The MOSAIC-SR information may be structured as a "minimum information checklist" conformant to standards such as those collectively gathered into the **MIBBI** recommendations; in this case MOSAIC-SR would be applied by implementing object models instantiating **MIBBI** policies. Alternatively, MOSAIC-SR reports can be derived from actual computer programs simulating research workflows, similar to **BIOCORDER**<sup>3</sup> (which is included by LTS, in an updated **C++** version, as one **CsAF** library). Finally, MOSAIC-SR presentations may be based on annotations applied to research/analytic code. For example, in the context of image analysis, the **PANDORE** project (an image-processing application) provides a suite of image-processing operations that can be called from computer code (together with the "Image Processing Objectives" Ontology mentioned earlier). Image-analysis pipelines can therefore be explained by annotating the pertinent function-calls (which **PANDORE** calls "operators") with terms from the **PANDORE** controlled vocabulary, providing annotation targets for MOSAIC-SR presentations.

MOSAIC-SR can express both computational workflows that are fully encapsulated by published code and real-world protocols concerning laboratory equipment and physical materials or samples under investigation. In the latter guise, MOSAIC-SR code can employ or instantiate standardized terminologies and data structures for describing experiments — such as **MIBBI** policies or **BIOCORDER** functions. In

<sup>2</sup>This toolkit provides a good case-study for research publication because it has an innovative **QT** and Common Workflow Language based extension mechanism; cf. [https://cbica.github.io/CaPTk/tr\\_integration.html](https://cbica.github.io/CaPTk/tr_integration.html).

<sup>3</sup>See <https://jbioleng.biomedcentral.com/articles/10.1186/1754-1611-4-13>.





this case, the role of **MOSAIC-SR** code is to serve as a serialization/deserialization endpoint for sharing research metadata. Conversely, when workflows are fully implemented within software developed as part of a body of research, **MOSAIC-SR** can provide a functional interface allowing this code to be embedded in scientific software. For the latter, **MOSAIC-SR** provides a framework for modeling how a software component specific to a given research project exposes its functionality to host and/or networked peer applications. There are also instances where both scenarios are relevant — the **MOSAIC-SR** code would simultaneously document real-world experimental protocols and construct a digital interface as part of a workflow which is part digital ("*in silico*") and part "real-world" ("in the lab").

### MOSAIC Annotations

Included in any **MOSAIC** plugin would be specially-designed **PDF** viewers for interactively reading authors' papers. In particular, these **PDF** applications would recognize cross-references pointing between publications and their associated supplemental archives. This cross-referencing would allow authors to identify concepts which are discussed and/or represented both in the research paper and in the archive, allowing them to annotate their papers more thoroughly. For example, the concept "**RNA** Extraction" may be discussed in a publication text, and also formally declared as one step in the lab processing as represented via **BIOCODER**, summarized in a **BIOCODER**-generated chart and included in the supplemental archive (a similar example is used in the documentation for **BIOCODER**). The **PDF** viewer would then ensure that the phrase "**RNA** Extraction" in the text is interactively linked to the concordant step in the experimental process, so that readers would then be able to view the **BIOCODER** summary as a context-menu action associated with the phrase, precisely where it appears in the **PDF** file (in other words, the annotation ground is a fine-grained location-designation in the text). Another example would be the phrase "Oxygenated Airflow", which refers to airflow in assisted-breathing devices, such as ventilators. To ensure that the device works properly, the equipment must be monitored to check that a steady stream of oxygen reaches the patient. Research into the design and manufacturing of ventilators and similar devices may then include "Oxygenated Airflow" both as a phrase within the article text itself and as a parameter in data sets evaluating the device's performance. In this situation, the publication-text location of the "Oxygenated Airflow" phrase should once again be annotated with links to the relevant part of the data set (e.g., a table column) where measurements of airflow levels are recorded.

In order to establish granular links between publication texts and data sets, **CsAF** introduces a novel "Hypergraph Text Encoding" (**HTXN**) system. This includes **L<sup>A</sup>T<sub>E</sub>X** code generators which output both **L<sup>A</sup>T<sub>E</sub>X** and **HTXN** representations of manuscripts, where the **L<sup>A</sup>T<sub>E</sub>X** files are provisioned with processing steps during **PDF** creation that generates geometric mappings between **PDF** viewport coordinates and **HTXN** annotations. Such information is then used by **CsAF** **PDF** viewers to cross-reference publications and data sets according to a protocol which allows users to conveniently switch back and forth between reading documents and visualizing data.

In addition to this technology for linking annotations with **PDF** and data-set annotations, **LTS** is working on frameworks for constructing annotations themselves with a rigorous semantics grounded on general-purpose models such as **OpenAnnotation** and the **Linguistic Annotation Framework**. We are, in particular, directing attention to bioimaging and other image-related use-cases via an "Annotation Exchange Format for Images" (**AXFi**), and to linguistic use-cases via an "Annotation Exchange Format for Text" (**AXFt**). The **AXFt** system supports graph-style annotations marking inter-textual connections on different linguistic scales. At the sentence level and smaller, these annotations are oriented to discourse-grounding models based on **Cognitive Grammar**; at the sentence level and larger, the primary model is **Sequence Package Analysis**, a discourse-representation methodology pioneered by **LTS**'s founder **Amy Neustein**. (More information about **HTXN** and annotations is provided in a separate **LTS** White Paper: <http://raw.githubusercontent.com/Mosaic-DigammaDB/CRCR/master/digamma/axf.pdf>.)



Systematic use of published scientific data requires accessing information on several different levels, including outlines of research methods and protocols as well as raw research data. Several standard models exist for how research data sets should be described, including Research Objects<sup>4</sup>, the **FAIR** (Findable, Accessible, Interoperable, Reusable) initiative<sup>5</sup>, **SciXML**<sup>6</sup>, and **SciData**<sup>7</sup>. Also relevant to scientific data sharing, though less general in scope, are formats for composing and/or annotating scientific publications, such as **IeXML**<sup>8</sup> and **JATS** (Journal Article Tag Suite). Despite (or perhaps because of) this diversity of representational possibilities, guidelines for constructing and annotating scientific data remain open-ended: researchers who want to rigorously annotate/document their published data can do so in many different ways, which presents challenges to software that interacts with research data.

Several different kinds of applications interface with scientific data, including software where such data is *created*; web services (and potentially **APIs**) where data is *hosted*; and programs which researchers use to find and download (or otherwise access) scientific data. Because data-representation paradigms vary, these applications need to be flexible enough to manipulate data and metadata structured in divergent ways. For example, each **MdsX** module is designed around a particular scientific repository/portal's protocols, so the structure of distinct **MdsX** modules may be noticeably different. Nevertheless, some general structuring principles can be developed with respect to querying and consuming repositories' resources. This motivates the idea of a "Scientific Data Repository Model" (**SDRM**), concretized within code implementing clients for individual portals (e.g., **MdsX** modules). The following section will present an overview of this **SDRM** representation, and subsequent sections will specify how **SDRM** is related to **ConceptsDB** and other technologies discussed here.

### Interoperating with Science Portals via the "Scientific Data Repository Model"

According to the **SDRM** representation, each data set is accompanied/constituted by a collection of "assets." These include raw data files; publication texts; research method descriptions; code for accessing/examining/analyzing research data; and annotations defining interconnections between these different elements. All material within or associated with a data set is classified into one of several **SDRM** "units," including those just enumerated and potentially other groupings.<sup>9</sup> **SDRM** units represent both serializations of data-set information and **GUI** components (such as Dialog Boxes) allowing readers of publications (and users of corresponding data sets) to view the unit's information. Moreover, **SDRM** units, by design, encapsulate their data in type-instances (e.g., objects of **C++** classes) that can be manipulated at runtime by computer code. In short, each unit has at least three representations (serial, **GUI**, runtime), giving rise to a spectrum of **SDRM** data that can be summarized schematically (see Table 1). In general, the data within **SDRM** units might be obtained from different sources, depending on the module — e.g., from publishers' **APIs**, or from downloaded data sets, or some combination thereof. (It is better to make at least part of the **SDRM** data available without having to download data sets as a whole, because a useful overview of the data set may contribute to whether users choose to download it *ab initio*.)

When **SDRM** information is provided as part of a dataset (rather than indirectly through an **API** or some other remotely-accessed resource), the **SDRM GUI** components can be integrated with the data

<sup>4</sup>See <https://eprints.soton.ac.uk/271587/1/research-objects-final.pdf>

<sup>5</sup>See [https://www.researchgate.net/publication/331775411\\_FAIRness\\_in\\_Biomedical\\_Data\\_Discovery](https://www.researchgate.net/publication/331775411_FAIRness_in_Biomedical_Data_Discovery).

<sup>6</sup>See <http://able.myspecies.info/scixml>

<sup>7</sup>See <https://stuchalk.github.io/scidata/>.

<sup>8</sup>See <https://www.semanticscholar.org/paper/IeXML%3A-towards-an-annotation-framework-for-semantic-Rebholz-Schuhmann-Kirsch/1d72a56b6576117c62f388a5f2193965e4c7e293>.

<sup>9</sup>For instance, data in some fields — such as genomics — sometimes becomes too large to be conveniently downloaded to an ordinary computer; in this case, portals such as GenBank offer query services to restrict very large data sets into small result-sets that researchers can examine directly. In these cases, the steps required to obtain information via large, remotely-hosted data steps would potentially be encapsulated into a distinct **SDRM** unit classified outside the main **SDRM** categories. Another potential example of **SDRM** units would be statements related to funding, potential conflicts of interest, and so forth.



SDRM Unit	Runtime Object	Serialization	GUI Component	LTS
Raw Data Files	File and archive objects identifying file types, preferred applications, software prerequisites, etc.	Structured review of archive layout (e.g., the Research Object Bundle <b>manifest</b> file)	Filesystem archive view plus information on file types and sizes	
Publication	Machine-Readable Text (e.g., <b>HTXN</b> )	Machine-Readable Serialization (e.g., <b>HTXN</b> -compatible <b>L<sup>A</sup>T<sub>E</sub>X</b> )	<b>PDF</b> Viewer	
Research Method	Research Protocol Encoding (e.g., <b>BIOCODER</b> )	Research Protocol Serialization (e.g., an <b>MIBBI</b> -based markup language)	Dialog Box to view "Minimum Information" checklist	
Dataset-Specific Code Library	Descriptions of library dependencies, compile steps, inter-type relations (e.g. correspondences between data objects and <b>GUI</b> types), etc.	Serialized build information	Dialog to set compile/user options; <b>IDE</b> integration	
Annotations on Data, Code, and Text	Integrated Annotation Object Model	Annotation Encoding (e.g. LTS's "Annotation Exchange Format")	Annotation-related context menus within <b>PDF</b> viewers and dataset <b>GUI</b> elements	

Table 1: Representations for SDRM Units

set’s overall front-end. In particular, if the data set provides its own customized **GUI** code — as is the case with Dataset Applications (discussed in the next section) — then such **GUIs** may include dialog boxes and similar components representing **SDRM** data.

Dataset Applications and the LTS Dataset Creator

Dataset Creator (**dsC**) is a framework for constructing data sets which include computer code based on the **QT** application-development platform. Dataset Creator allows scientists to publish research data in conformance with standards such as Research Objects, while also equipping those using the data with code libraries and interactive software, to stimulate adoption and reuse of the data (and the methods or perspectives which inform it). **QT**, the leading native cross-platform development toolkit, is a comprehensive framework encompassing a thorough inventory of programming features — networking, **GUI** implementation, file management, data visualization, **3D** graphics, and so forth. Data sets based on **QT** require users to obtain a copy of the **QT** platform, but **QT** is free for non-commercial use and easy to install — importantly, **QT** is wholly contained in its own folder and does not affect any other files on the user’s computers (in this manner **QT** is different than most software packages, which usually demand a “system install”).

By leveraging the **QT** platform, **dsC** enables standalone, self-contained, and full-featured native/desktop applications to be uniquely implemented for each data set, distributed in source-code fashion along with raw research data. Adopting such a data-curation method makes data sets easier to use across a wide range of scientific disciplines, because the data sets are freed from having to rely on domain-specific software (software which may be commonly used in one scientific field but is unfamiliar outside that field). Ideally, dataset code may be distributed both via plugins to domain-specific software popular for a given research field, and also via standalone Dataset Applications accessible to researchers coming from other disciplines.

Because every data set is unique, each Dataset Application will necessarily include some code specific to that one Research Object. However, **dsC** will provide a core code base and file layout which is shared by all **dsC** data sets by default. This common core is structured in part by the goal of developing Dataset Applications in a **QT** context; for instance, **dsC** projects are structured to use **QT**’s “**qmake**” build system as the primary tool for compiling data-set code. The common **dsC** code therefore includes **qmake** project files which support compiling application with several build configurations. In this framework, data-set users are classified into several different roles — in addition to ordinary users (specifically, researchers who want to work with and draw information from data sets but have no development connection to these data sets themselves), **dsC** recognizes roles for authors, editors, testers, and other users who are responsible for bringing data sets into publication-ready form to begin with. Depending on the administrative role, data set code can be compiled with additional features (e.g., unit testing features).



Another core component of **dsC** is **L<sup>A</sup>T<sub>E</sub>X** code that authors may use when preparing documents accompanying their data set. These **L<sup>A</sup>T<sub>E</sub>X** files encompass special functionality for defining code annotations and semantically significant points in article text, such as sentence and paragraph boundaries. This **L<sup>A</sup>T<sub>E</sub>X** code can be used in conjunction with a pre-processor that generates **L<sup>A</sup>T<sub>E</sub>X** files from a special input language. The goal of these text-processing technologies is to improve the interoperability between research papers and data sets. In particular, the **L<sup>A</sup>T<sub>E</sub>X** pre-processors (and subsequent **L<sup>A</sup>T<sub>E</sub>X**-to-**PDF** converters) generate files in a custom **CsAF** hypergraph-serialization format (**HGXF**, for "Hypergraph Exchange Format,") which store information about textual and **PDF** viewport coordinates specifying the location of semantically meaningful elements such as sentences and annotations. These files are then zipped and embedded in **PDF** files. With **dsC**, the resulting **PDF** files can then be loaded into a customized **PDF** viewer capable of reading the embedded **HGXF** data. This allows the **PDF** application to utilize the embedded information so as to provide a more interactive reading experience — for instance, viewing annotations or copying sentences via context menus, where viewport data maps cursor position to textual elements visible on the **PDF** page. These features provide an application-level interface between the **PDF** viewers (considered as **GUI** components) and the corresponding **GUI** components in Dataset Applications.

With proper customization, both the **PDF** viewers and the **dsC** Dataset Applications can interoperate, with **PDF** context menus calling up windows in the Dataset Application's **GUI** implementation, and vice-versa. For example, researchers reading the **PDF** version of a scientific article can launch the Dataset Application to explore some detail mentioned in the text. This is an example of where micro-citations are practically useful: any microcitable element in a document (such as a table, column, row/record, or analytic procedure formalized as a procedural asset associated with a data set) can be linked to a corresponding **GUI** element in the Dataset Application. For example, a statistical parameter — mentioned by name in the text, and perhaps represented in serialization within raw data — can be mapped to a **GUI** table column, and specifically the column header; this is then an annotation target, in the sense that for readers to gain more information it is proper to link mentions of the relevant scientific concept in article text to the column header as a graphical element that can be made visible. The link is operationalized by implementing a procedure to show the **GUI** window where the table is located, and ensure that the column header lies in the visible portion of the screen, as a response to readers on the **PDF** side signaling a desire for information on the annotated text element. In the opposite direction, database elements can be annotated with links that can be used by the Dataset Application to launch a **PDF** window opened to the page and location where the corresponding concept is discussed in the article.

In order to properly model this semantic, viewport, and data set data integration, **dsC** uses **HTXN**, described earlier, for document representation. With **dsC**, **HTXN** files are not only associated with data set assets; they are also machine-readable document encodings that can be introduced into publication repositories and other corpora oriented toward machine-readability. Authors can host **HTXN** files within data sets and link to them via services such as CrossRef, thereby ensuring that a highly structured, machine-readable version of their papers is available for text and data mining. The **HTXN** protocol is also useful for encoding natural-language content which becomes part of a data set as data assets in themselves; for example, linguistic case-studies, or (in the biomedical context) patient narratives. To establish annotation targets for cross-references between publications and data-set code, **dsC** employs a code-annotation system based on hypergraphs, along with (on the manuscript side) the Annotation Exchange Format (**AXF**) cited above (see Page 5).

### Hypergraph Data Serialization and "Grounding"/"Meta-Serialization"

As a database engine, **ConceptsDB** is focused on data persistence — that is, storing information so that it may be shared among different instances of one or several applications; and preserving information between application-sessions, where each "session" refers to the time during which one user launches, uses, and ultimately exits an application. In addition to data persistence, however, **ConceptsDB** also encapsulates theoretical research into optimal strategies for data *modeling*, which





includes representations during application runtimes and serializations for sharing data between heterogeneous end-points. **ConceptsDB** introduces the idea of “infoaset classes,” which are intermediate structurations of data type-instances, organized so that interconnected clusters of infoaset class objects can be modeled as hypergraphs. Infoaset classes are guidelines for marshaling runtime values into structures that can be persisted in **ConceptsDB**. At the same time, infoaset classes can also be used as preliminary encodings for data sharing and serialization.

The **ConceptsDB** Application Framework (**CsAF**) is designed around the premise that the primary goal of any data modeling system (and, as a surface-level encoding of data to be shared, any markup language) is to encapsulate data which is used by and present in one application, to be subsequently re-used by a different application — or by the same application at a future time. In any data-modeling scenario there is an application which *creates* a data structure, and an application (not necessarily a different one) that *consumes* that data structure. “Consuming” in this context means reading the data so that it may be incorporated into some computational process and/or rendering the data so that it may be viewed by a human user. Of course, any consumer-application may double as a creator, adding new data or revising that which it has consumed.

Data sharing is “correct” insofar as the consumer-application reads, understands, displays, and allows human users (as well as potentially computer algorithms) to manipulate the data in ways which are functionally equivalent to the creator-application. When this functional equivalence is achieved, the two applications can be said to be *aligned*. Creator-to-consumer alignment can happen in one of two ways. One option is that the data model shared between the two applications is formalized so that different applications are naturally interoperable because they are implemented according to similar guidelines. This may be called “coincidental alignment.” There are, in turn, two pathways toward coincidental alignment. On the one hand, data-model formalization may stipulate criteria on creator and consumer applications to ensure that any creator is aligned with any consumer. Such a data model would essentially provide a “checklist” of requirements to be satisfied, procedures to be implemented, and so forth, so that applications are only considered to be creators and/or consumers by fulfilling the checklist contract.

On the other hand, protocol formalization may focus on the medium through which data shared between creator and consumer is structured and encoded. This second mode of standardization — coincidental alignment based on data-serialization formats — is arguably the most common strategy for enforcing alignment (even if it is not identified in these terms). For example, **XML** schemata represent constraints on any creator of information encoded by conformant **XML** documents. Given these constraints, **XML** consumers can guarantee alignment with creators by incorporating the schema as a contract the creator obeys. The consumers need not have any knowledge of the creators *apart from* the fact the **XML** they create conforms to that applicable schema.

Whether by explicit application-level requirements or via contracts on the medium by which shared data is encoded (e.g., **XML** schemata), coincidental alignment can occur without any explicit coordination between the code used by creator and consumer applications, respectively. As an alternative model, however, we can identify *deliberate* alignment, wherein the consumer aligns with a creator by explicitly referring to and mimicking the creator’s data models and procedures. To facilitate deliberate alignment in this sense, creator applications can architect the relevant code to prioritize transparency and reuse — e.g., factoring this code into a semi-autonomous library which the consumer application can use directly, or use as the basis of its own library.

In contrast to many standard data meta-models, **CsAF** openly embraces *deliberate* rather than coincidental alignment in this sense. Data-sharing protocols are not only necessary application features; they are also proxies for a scientific/technical formalization of the principles underlying a given data profile. Standardizing only a surface-level realization (e.g., **XML** schemata) fails to rigorously capture this scientific dimension. One way to document the shape of any relevant data is simply to explain how an **XML** document will be structured insofar as it encodes data accordingly<sup>10</sup>; however, the structure of an **XML** document does not, in and of itself, present a clear picture of how the infor-

mation which the document represents is semantically organized. Even though **XML** is processed by computer programs, one cannot even determine from an **XML** document or schema which **XML** elements (if any) correspond to data types recognized by applications which read and/or write the corresponding **XML** code.<sup>11</sup> In short, because applications can read or use data from an **XML** document in different ways, the document's structure does not itself provide a clear picture of how the specific code that reads the **XML** is organized. Such uncertainty becomes significant when one wishes to use markup specifications as an indirect strategy for documenting parameters and features of the data structures which are serialized via the relevant markup language. In effect, markup/text-based serialization operates on two levels: on the one hand, the specific document provides an encoding of data conforming to a given structure; but, at a more abstract level, one can model the relationship between the surface-level document structure and the application-level data structures that are thereby serialized. This second level of detail is usually implicit and unstated, but in **CsAF** technology, such "meta-serialization" — a term we are using to suggest the idea of providing meta-data *about* a serialization — is directly formulated through a notion of "grounding," which involves adding supplemental markup clarifying how documents instantiating a serialization format relate to the coding protocols and data types of software that reads or creates these documents.

For maximum expressivity, **CsAF** introduces a meta-serialization system that can be applied to languages more flexible than **XML**, notably **TAGML**.<sup>12</sup> In effect, **CsAF** will implement parsers for an extended version of **TAGML**, one that includes an additional "grounding" layer. Grounding, in this context, means describing how elements in the markup — nodes, attributes, and character sequences — are "grounded" in application-level types, data fields, and other programming constructs. **CsAF** provides parsers for this "Grounded **TAGML**" (**G<sub>TAGML</sub>**) language as well as converters to output serialized data in more conventional formats (e.g., **XML**). **CsAF** will also provide code libraries for extracting data from **G<sub>TAGML</sub>** documents.

By design, **G<sub>TAGML</sub>** schemata operate on two levels: first, they constrain the organization of **G<sub>TAGML</sub>** files themselves; and second, they stipulate how **G<sub>TAGML</sub>** document structure relates to the type systems of code libraries serializing and deserializing **G<sub>TAGML</sub>** files. To model the second level of metadata, **G<sub>TAGML</sub>** introduces a hypergraph-based type model organized around "infosets," which are structured overviews of application-level data types. To fully utilize **G<sub>TAGML</sub>** features, programmers may consequently compose info-set classes (described earlier), as wrappers around ordinary data types (e.g., **C++** classes), whose instances are serialized via nodes or character strings within **G<sub>TAGML</sub>** documents. Info-set classes provide hypergraph "views" onto type-instances, and act as a bridge between data types and their associated **G<sub>TAGML</sub>** schemata.

Because **G<sub>TAGML</sub>** is not only a hypergraph-serialization format, but also a representational paradigm for "meta-serializaton" data, **G<sub>TAGML</sub>** is closely interconnected with **QDB** data models. **QDB** recognizes numerous forms of hypergraph structuring elements and levels of organization; each of these can be annotated via **G<sub>TAGML</sub>** to document the relation between **G<sub>TAGML</sub>** node patterns and **ConceptsDB** hypergraph formations (or any construction which can be mapped to **QDB**). More information on the structural principles underlying **QDB** hypergraph categories can be found in an LTS proposal for applying **QDB** in bioimaging contexts (see <http://raw.githubusercontent.com/Mosaic-DigammaDB/CRCR/master/hmcl/proposal.pdf>).

*For more information please contact:*

**Amy Neustein, Ph.D., Founder and CEO**

**Linguistic Technology Systems**

**amy.neustein@verizon.net • (917) 817-2184**

<sup>10</sup>Potentially, such specification can be a single **XML DTD** file, or an **XML** sample, providing a convenient reference point for developers to grasp the underlying data model.

<sup>11</sup>For example, the **XML** portion of **OME-TIFF** (the principal Open Microscopy Environment imaging format) includes an explicit **Image** element (which gathers up all significant image metadata); an application reading **OME-TIFF** files might, therefore, introduce a single datatype — analogous to (and maybe wrapping an) **itk::Image** from the Insight Toolkit imaging libraries — bundling the data in that part of the **OME XML**. In this case, there will be a one-to-one correspondence between **XML** structure and application-level data types, at least for that one **Image** node. On the other hand, software reading **OME-TIFF** information might not manipulate images directly, but rather pull out other kinds of metadata, such as an experimenter's name or description of the microscopic setup. In this case, the application might not have an explicit "object" representing the image itself, but it may still read information about the image from child nodes of the **XML Image** element.

<sup>12</sup>See <http://www.balisage.net/Proceedings/vol21/print/HaentjensDekker01/BalisageVol21-HaentjensDekker01.html>.

