

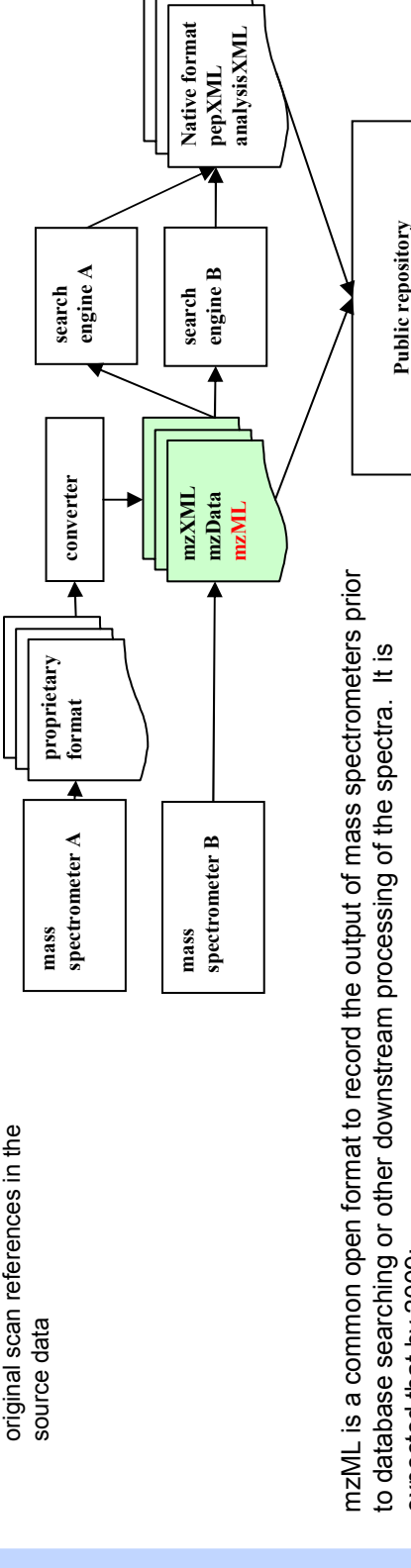
Design and implementations of the new HUPO Proteomics Standards Initiative's mass spectrometer output file standard format: mzML 1.0

Eric W Deutsch¹, Pierre-Alain Binz², Darren Kessner³, Matt Chambers⁴, Luisa Montecchi-Palazzi⁵, Jim Shofstahl⁶, Josh Tasman¹, Randall K Julian⁷, Fredrik Levanders⁸, Puneet Souda⁹, and Lennart Martens⁹
¹Institute for Systems Biology, Seattle, WA; ²Swiss Institute for Bioinformatics and Geneva Bioinformatics, Geneva, Switzerland; ³Cedars-Sinai Medical Center, Los Angeles, CA; ⁴Vanderbilt University, Nashville, TN; ⁵European Bioinformatics Institute, Hinxton, UK; ⁶Thermo Fisher, San Jose, CA; ⁷Indigo Biosystems, Carmel, IN; ⁸Lund University, Lund, Sweden; ⁹University of California Los Angeles, Los Angeles, CA

Overview

mzML is a new data format for the storage and exchange of mass spectrometer output files. It follows on the successful mzXML and mzData formats. mzML has been designed by merging the best aspects of both previous formats into a single unified format that is intended to replace all earlier formats.

- Version 1.0.0 just released
- Accompanied by a controlled vocabulary and semantic validation rules
- Many implementations of the format already exists, insuring quick adoption of the format
- Developed with full participation of academic researchers, hardware and software vendors
- Vendors have committed to supporting the new format once released.
- Format has been tested with several instance documents and minor implementations of the format during beta testing
- mzML is expected to replace mzXML, and mzData, but not expected to completely replace vendor binary formats.
- NativeID references back to the original references in file source data

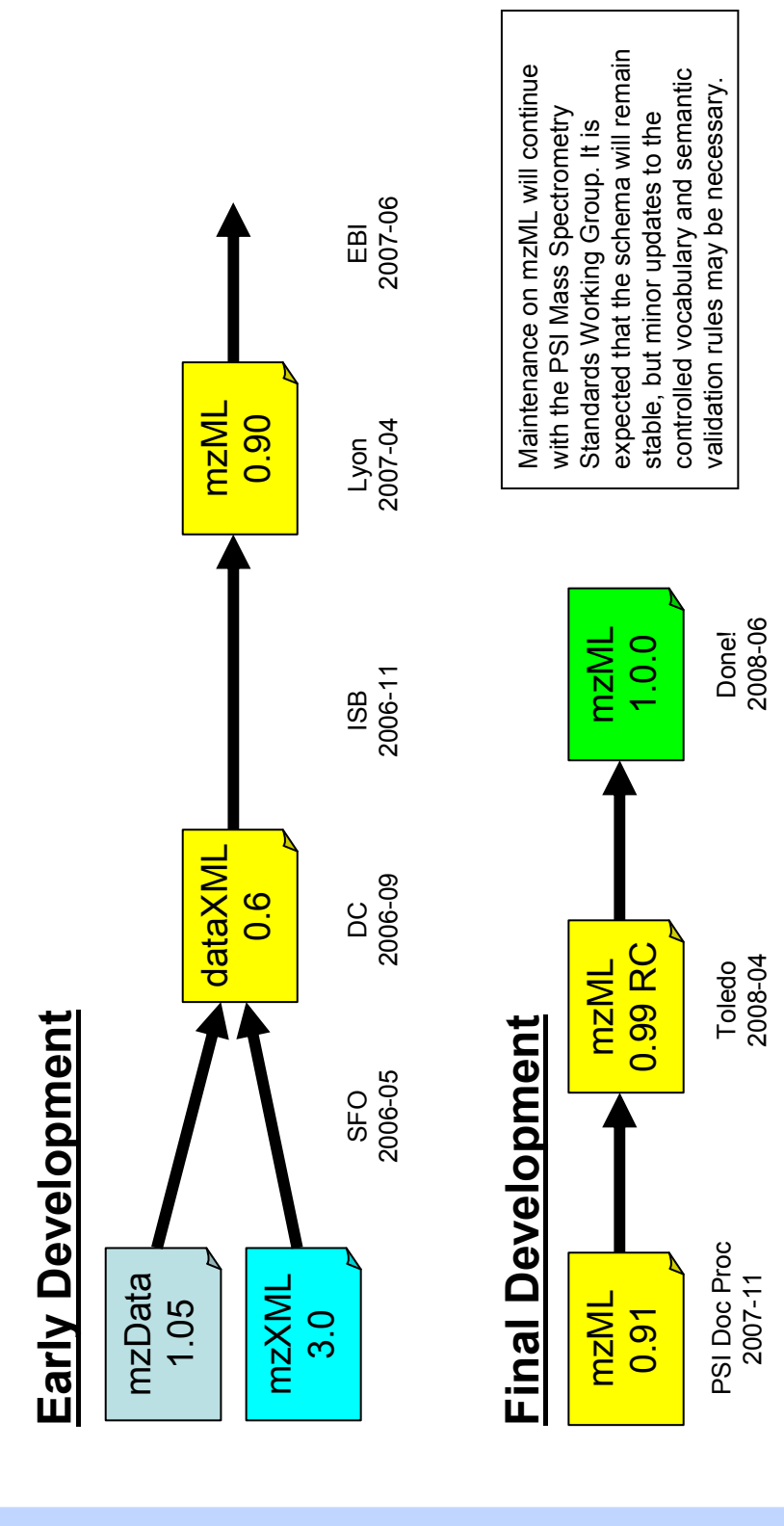


mzML is a common open format to record the output of mass spectrometers prior to database searching or other downstream processing of the spectra. It is expected that by 2009:

- Instrument vendors will write out or convert to mzML
- Search engines or other spectrum processing software will read and process mzML
- Data repositories will accept, process, and store mzML documents

History

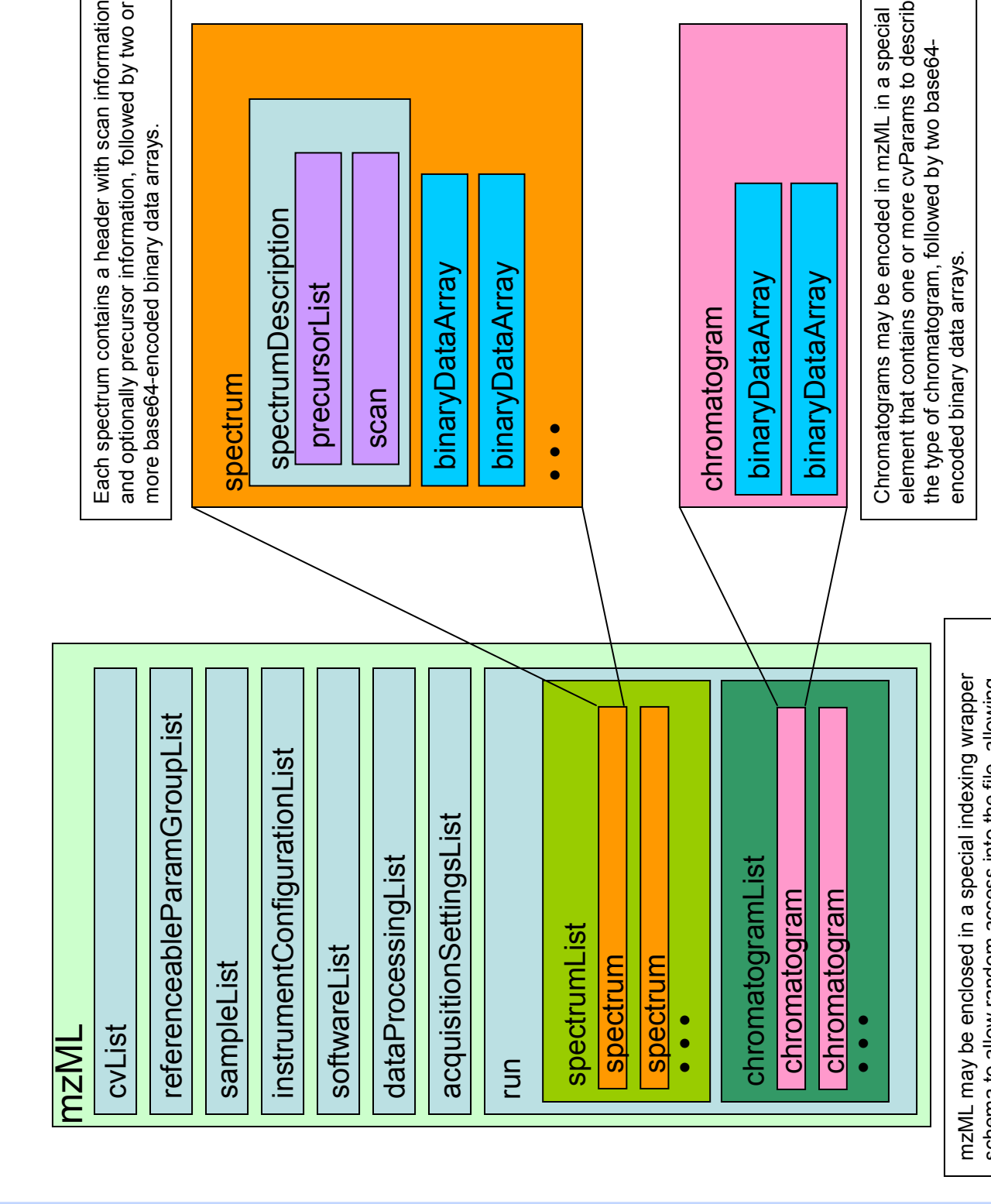
mzML has been under development for two years with full participation of academic researchers, hardware and software vendors, mass spectrometer manufacturers, and search engines. The format was first developed in 2007, and was submitted to the HUPO Proteomics Standards Initiative in 2007. The format was submitted to the HUPO Proteomics Standards Initiative in 2007. The format was submitted to the HUPO Proteomics Standards Initiative in 2007. The format was submitted to the HUPO Proteomics Standards Initiative in 2007.



Meanings on mzML will continue with the PSI Proteomics Standards Working Group. It is expected that the schema will remain stable, but minor updates to the controlled vocabulary and semantic validation rules may be necessary.

Schema Outline

The mzML schema is designed to contain all the information for a single MS run, including meta data about the spectra plus all the spectra themselves (raw or processed). The reader of the file or the file archives information about the source of the data as well as information about the sample, instrument and software that processed the data.



Chromatograms may be encoded in mzML in a special element that contains one or more cvParams to describe the type of chromatogram, followed by two base64-encoded binary data arrays.

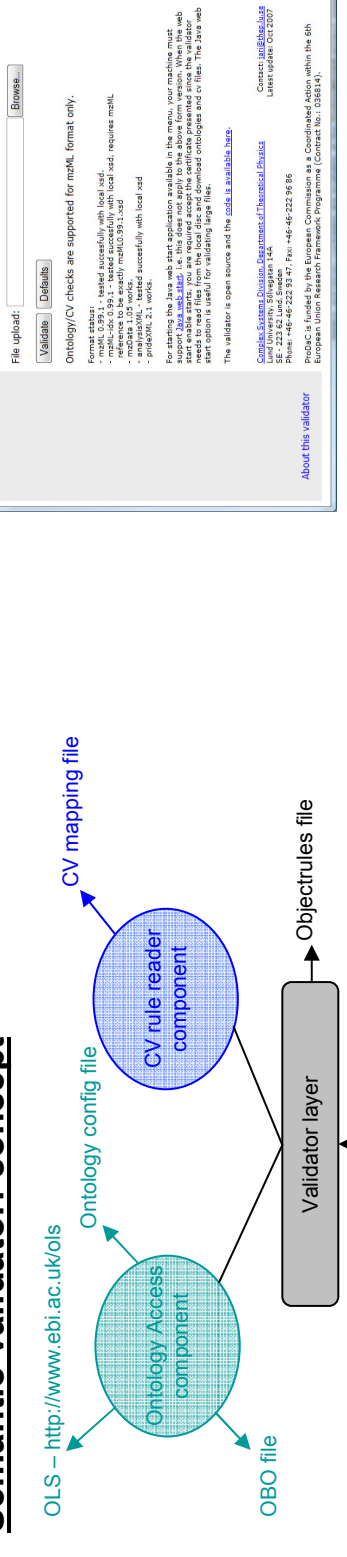
mzML may be enclosed in a special indexing wrapper schema to allow random access into the file, allowing software to pull out one or more arbitrary spectra.

Semantic Validator

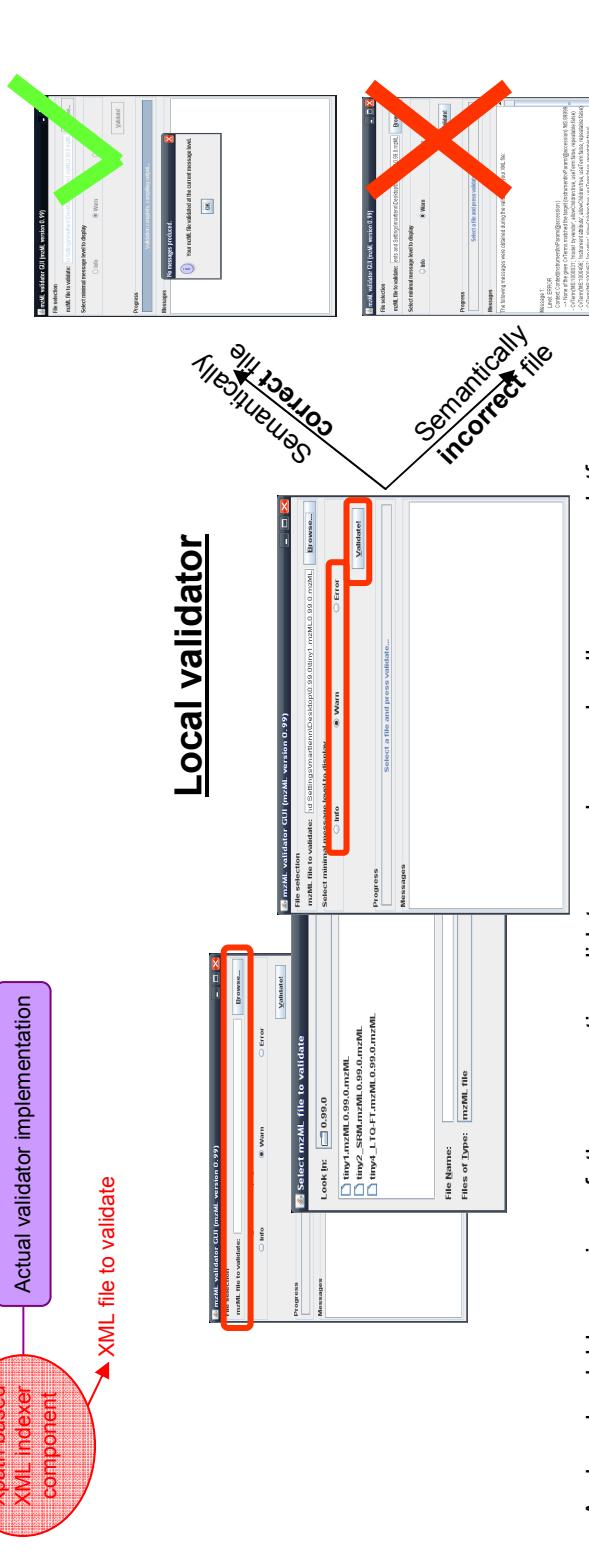
One of the benefits of the previous mzData format was its considerable flexibility in allowing writers of the format to encode additional information (e.g. vendor specific information) that was not covered by the standard. However, this flexibility was not maintained in the mzML format. The semantic validator was developed to address this problem with the semantic validator tool for mzML. The semantic validator can:

- Ensure that an mzML document is well formed and conforms to the mzML schema
- Ensure that controlled vocabulary terms are used in the correct places in the document
- Allow alternate rules based on the type of data being written
- Allow alternate levels of compliance (e.g. basic mzML, MIAPE-MS compliant mzML)
- Semantic rules can be updated along with the controlled vocabulary without changing the schema
- Available as a web page (see below) or as a standalone tool

Semantic validator: concept



Online validator at the ProteoWizard site allows anyone to perform semantic validation on any mzML file



A downloadable version of the semantic validator can be run locally on any platform supporting Java to validate files without needing to transmit them to a remote web site.

Example Instance Documents

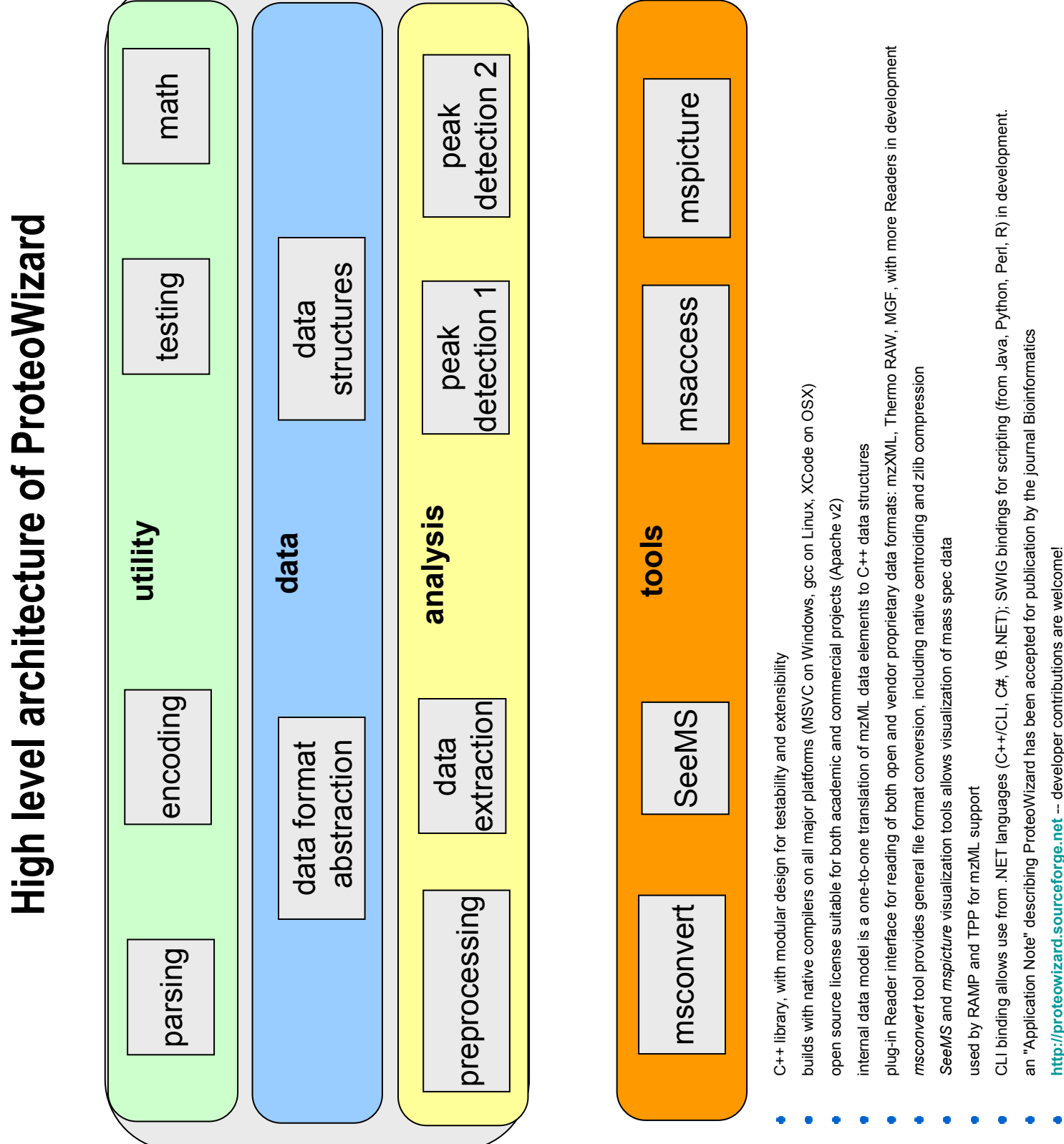
In order to exercise the schema and demonstrate that the various use cases have been adequately modeled, we have developed several example instance documents. Some of the documents are hand-crafted with an ordinary editor, while others are written out as a software test as part of the ProteoWizard reference implementation. In addition, several instance documents are conversions of real data files using the reference or other implementations of converters.

A sample snippet of an example mzML document, showing the top header portion of the file.

ProteoWizard

The ProteoWizard software project, initiated by the Spielberg Family Center for Applied Proteomics at the Cedars-Sinai Medical Center, is a multi-platform open source project that provides a comprehensive set of tools for the processing and analysis of mass spectrometry data. The tools are designed to be easy to use and to integrate with existing software. The ProteoWizard software project provides a comprehensive set of tools for the processing and analysis of mass spectrometry data. The tools are designed to be easy to use and to integrate with existing software.

High level architecture of ProteoWizard



- C++ library with modular design for testability and extensibility
- builds with native compilers on all major platforms (MSVC on Windows, gcc on Linux, XCode on OSX)
- open source license suitable for both academic and commercial projects (Apache v2)
- internal data model is a one-to-one translation of mzML data elements to C++ data structures
- plug-in Reader interface for reading of both open and vendor proprietary data formats: mzXML, Thermo RAW, MGF, with more Readers in development
- recover tool provides general file format conversion, including native controlling and zlib compression
- SeeMS and mspicture visualization tools allow visualization of mass spec data
- used by RAAMP and TPP for mzML support
- CLI handling allows use from .NET languages (C++/CLI, C#, VB.NET), SWIG bindings for scripting (from Java, Python, Perl, R) in development.
- an "Application Note" describing ProteoWizard has been accepted for publication by the journal Bioinformatics
- <http://proteowizard.sourceforge.net> - developer contributions are welcome!

Documentation

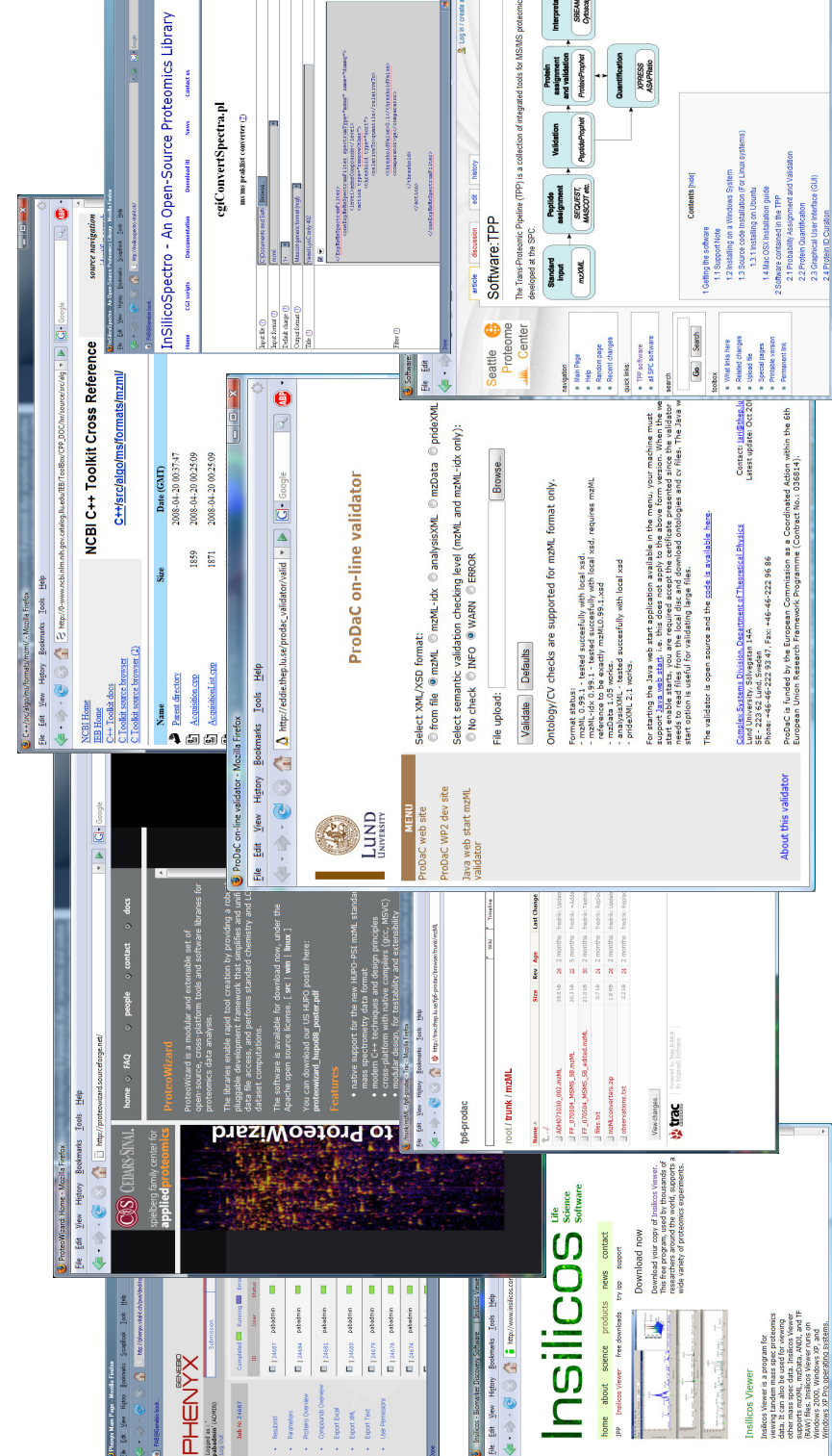
The full specification of the format is presented in a specification document that describes various aspects of the format as well as all details of the format elements. HTML and PDF documentation are generated by programmatically combining 4 different components: 'xsd' schema, controlled vocabulary, semantic validation mapping file, example documents.

This page of the full mzML Specification Document

Software Implementations

The best way to test a new format is by implementing it in software. Ideally, as a format is implemented, one finds minor inconsistencies or missing features. The initial release of mzML is strengthened by the breadth of implementations that already exist and have exercised the various use cases:

- ProteoWizard C++ library reference implementation: reads mzML, and writes mzML. Can convert mzXML and RAW to mzML.
- RAAMP (Random Access Minimal Parser) C library can read mzML, mzXML, mzData files via same API. mzML reading performed with ProteoWizard
- Trans Proteomic Pipeline (TPP) can read and process data in mzML format via RAAMP
- ISI format converters: RoADW (Thermo), Wof (Waters), mzWf (ABI/MDS), Trapper (Agilent)
- Thermo Fisher beta RAW -> mzML converter
- Semantic Validator and Java library reads and validates that a document is semantically correct
- Phenix search engine can read and search spectra in mzML format
- NCBI C++ mzML reader classes
- Insilicos Viewer - file browser and spectrum viewer can read and display spectra from mzML, mzXML, mzData, RAW formats
- SeeMS file browser, spectrum viewer, chromatogram viewer, annotator for mzML, mzXML
- Proteo Software Environment includes converters for peak lists of various formats to mzML, and performs reading of mzML files
- InsilicoSpectro open source library (Perl) has a spectrum file format conversion tool that reads mzML
- Hermes mzML -> mzData -> mzXML converter (Java)



Conclusions

The mzML format is now complete and mzML 1.0.0 is released. We encourage all authors and vendors to begin supporting this new format in new and updated software. The format includes the best features from pre-existing open formats and has additional support for chromatograms and some other features deemed highly desirable.

It is expected that the schema will remain stable for at least a year, hopefully more. However, the controlled vocabulary and semantic validation rules will continue to be updated and refined as all authors and vendors finish implementing their software for mzML.

To learn more, see the mzML Development Page:

<http://psidev.info/index.php?q=node/257>

The mzML effort has involved many people in the PSI and in the community. We gratefully acknowledge the contributions of:

- Jari Hakkinen (Lund)
- Jason Falkner (U Michigan)
- David Horn (Agilent)
- Brian Pratt (Insilicos)
- Phil Jones (Cardiff)
- Ruth McNally (Cardiff)
- Ron Beavis (JBC)
- Norman Paton (Manchester)
- Ruedi Aebersold (ETHZ)
- Mark Stum (U Tuebingen)
- Parag Mallick (CSHS)
- Chris Taylor (EBI)
- Patrick Pedrioli (ETHZ)
- Sean Seymour (ABI)
- Rune Philsof
- David Cressy (Matrix Science)
- Wilfred Tang (ABI)
- Howard Reed (Waters)
- Jim Langridge (Waters)
- Manus Kallhaer (Bruker)
- PSI Steering Group
- PSI Participants