• Structural • Informatics Group

Introduction

In clinical research studies, multiple applications (i.e. data sources), each with a specialized function, are often used to track information about patients, their specimen samples and lab results. Simple questions regarding a patient's records can be laborious to answer. A number of bioinformatics integration platforms currently exist that can facilitate the answering of such questions (Amalga, i2b2) but are often costly to install, configure and maintain.

We present a reusable, lightweight, framework and workflow that facilitates the continuous integration and querying of study data from heterogeneous sources in order to answer routine, operational questions. The goal of our project is to develop an inexpensive and flexible approach to cross-application data integration that will provide informaticists with the tools they need to assist researchers. Biomediator Data Integration System (BDIS) serves as the reference implementation for our work.

2 ml second an U1 AU an as din a UUTE AU2	2
xml version="1.0" encoding="UTF-8"?	xml version="1.0" encoding="UTF-8"?
<pre><xs:schema <="" attributeformdefault="unqualified" elementformdefault="quali" pre=""></xs:schema></pre>	<pre><xs:schema <="" attributeformdefault="unqualified" elementformdefault="quali" pre=""></xs:schema></pre>
<pre><xs:element name="families" type="familiesType"></xs:element></pre>	<xs:element name="records" type="recordsType"></xs:element>
<pre><xs:complextype name="familiesType"> [4 lines]</xs:complextype></pre>	<xs:complextype name="itemType"></xs:complextype>
<xs:complextype name="subject-idType"> [6 lines]</xs:complextype>	<xs:sequence></xs:sequence>
<pre><xs:complextype name="aliquot-used-byType"> [6 lines]</xs:complextype></pre>	<pre><xs:element name="study" type="xs:string"></xs:element></pre>
<xs:complextype name="sample-idType"> [6 lines]</xs:complextype>	<pre><xs:element name="redcap_event_name"> [12 lines]</xs:element></pre>
<pre><xs:complextype name="aliquotType"> [67 lines]</xs:complextype></pre>	<pre><xs:element name="altid" type="xs:string"></xs:element></pre>
<pre><xs:complextype name="cell-countType"> [6 lines]</xs:complextype></pre>	<pre><xs:element name="group"> [11 lines]</xs:element></pre>
<pre><xs:complextype name="sampleType"> [25 lines]</xs:complextype></pre>	<pre><xs:element name="entrysubcat"> [11 lines]</xs:element></pre>
<pre><xs:complextype name="aliquot-used-dateType"> [6 lines]</xs:complextype></pre>	<pre><xs:element name="slesubcat"> [10 lines]</xs:element></pre>
<pre><xs:complextype name="aliquot-used-forType"> [6 lines]</xs:complextype></pre>	<pre><xs:element name="sscsubcat"> [11 lines]</xs:element></pre>
<pre><xs:complextype name="subjectType"> [50 lines]</xs:complextype></pre>	<pre><xs:element name="enroll_date" type="xs:string"></xs:element></pre>
<pre><xs:complextype name="familyType"> [9 lines]</xs:complextype></pre>	<pre><xs:element name="dob" type="xs:string"></xs:element></pre>
<pre><xs:complextype name="colType"> [6 lines]</xs:complextype></pre>	<pre><xs:element name="ethnicity"> [9 lines]</xs:element></pre>
<pre><xs:complextype name="rowType"> [6 lines]</xs:complextype></pre>	<pre><xs:element name="entryrace"> [14 lines]</xs:element></pre>
<pre><xs:complextype name="legacy-codeType"> [6 lines]</xs:complextype></pre>	<pre><xs:element name="enroll_comm" type="xs:string"></xs:element></pre>
<pre><xs:complextype name="x-refType"> [6 lines]</xs:complextype></pre>	<pre><xs:element name="enroll_subject_tracking_information_complete"></xs:element></pre>
<pre><xs:complextype name="relationship-codeType"> [6 lines]</xs:complextype></pre>	<pre><xs:element name="qsex"> [8 lines]</xs:element></pre>
<pre><xs:complextype name="created-atType"> [7 lines]</xs:complextype></pre>	<pre><xs:element name="agesympt" type="xs:string"></xs:element></pre>
<pre><xs:complextype name="idType"> [6 lines]</xs:complextype></pre>	<pre><xs:element name="raynauds"> [8 lines]</xs:element></pre>
<pre><xs:complextype name="other-codeType"> [6 lines]</xs:complextype></pre>	<pre><xs:element name="raydate" type="xs:string"></xs:element></pre>
<pre><xs:complextype name="updated-atType"> [7 lines]</xs:complextype></pre>	<pre><xs:element name="skintight"> [8 lines]</xs:element></pre>
<pre><xs:complextype name="DOBType"> [7 lines]</xs:complextype></pre>	<pre><xs:element name="skindate" type="xs:string"></xs:element></pre>
<xs:complextype name="age-at-onsetType"></xs:complextype>	<pre><xs:element name="rash"> [8 lines]</xs:element></pre>
Fig. 1A.) XML schema for REDCap.	Fig. 1B.) XML schema for rLab.
r 1g. 17 j Mine Schema for KEDCap.	1 15. ID. J MILL SCHEIMA IOI ILAU.
- 2xml vencion "1 Q" encoding "UTE 9"2	



Lightweight Data Integration Frameworks for Clinical Research Ronald Shaker, BS¹, Xenia Hertzenberg, BS¹, James F. Brinkley, MD, PhD^{1,2,3} ¹Departments of Biological Structure, ²Computer Science and Engineering, ³Medical Education and Biomedical Informatics, University of Washington, Seattle, WA

Methods

We begin by narrowing requirements for the applications we wish to integrate to a reasonable scope. By keeping expectations simple, data sources not meeting our basic criteria may more easily be altered to provide access to their data.

Queries and Semantic Mappings:

BDIS uses a mediated¹ approach to integrating heterogeneous data. A common schema, defined as an XSD, unions the elements of interest from each data source (Fig 1). An XQuery is then used to simultaneously filter and map results onto the mediated schema (Fig 2). XML is used as the output syntax of all data sources. Each source must formally define its output schema in the form of an XSD. Schema definition promotes semantic understanding of the data source, provides validation of input data sets, and can also facilitate the automated generation of user interfaces for filtering and constraining query results.

Data Acquisition:

Sources are required to provide an API for downloading their data as XML. Connectivity to each source is established by specifying a network address, HTTP method, authentication credentials and query parameters. Recurring extractions can be scheduled at regular intervals to provide near real-time reporting.

Reporting and Security:

XSL transformations may be associated with each query and are automatically applied to XML query results, generating user-friendly formats such as CSV and HTML. Original XML results, as well as transformations, are made available to external clients via a RESTful API

Externally, access is controlled using user accounts and single-access tokens over HTTPS. Internally, each piece of sensitive data is assigned to an owner and group. Owners may share their data by adding users to the appropriate group. User access to data is controlled at the ORM level by use of an ActiveRecord plug-in in Ruby and an added security layer above JPA in Java.

- <spreadsheet>

<subject-code>JRA140A</subject-code>

<sample-code>JRA140A</sample-code>

<sample-type>BLOOD</sample-type>

<aliquot-type>PLASMA</aliquot-type

<entrysubcat>RF+ Poly</entrysubcat>

<ethnicity>Non Hispanic/Latino</ethnicity:

<sampletype___0>0</sampletype___0>

<sampletype___1>0</sampletype___1>

<sampletype___2>0</sampletype___2>

<redcap_event_name>Initial Visit</redcap_event_name>

<enroll_subject_tracking_information_complete>Complete<

<sample-comments/>

<state>NOT_USED</state>

<aliquot-comments/>

<study>JRA140</study>

ltid>RRR009</altid>

<entryrace>White</entryrace>

group>JIA</group>

<sscsubcat/>

<enroll comm/>

<sample_date/>

<age_at_donation/>

<cell-count/>

<col>1</col>

<row>I</row>

<box>217</box>

- <row>

REDCap	ITHS Institute of Translationa		leh C	sion						
Logged in as rshaker@washington.edu Log out	Live V institute of translational Health Sciences									
	UW Seattle									
My Projects	Institute of Translational Health Science (ITHS) Stevens 2107: Family Study of Pediatric Autoimmunity									
Project Home Project Setup Project status: Production										
Data Collection	📧 Data Entry: Event Grid									
Data Entry	The grid below displays the form-by-form progress of da You may click on the colored buttons to access that form				-				•	
Applications	to the Define My Events page.									
Data Export Tool Data Import Tool Data Import Tool Data Comparison Tool	Study ID JRA021									
File Repository	Events									
Graphical Data View & Stats API Report Builder	Data Collection Instrument	Initial Visit (1)	Visit 2 (2)	Visit 3 (3)	Visit 4 (4)	Visit 5 (5)	Visit 6 (6)	Visit 7 (7)		
Remente	Enroll Subject - Tracking Information	۲								
Reports	Sample Information	۲	۲	۲	٠	۲	۲			
1) Novo Interests 2) Selections	Pediatric Autoimmunity Study Questionnaire	۲								
3) SLEDAI PDL1	Pediatric Autoimmunity Study Follow-up Questionnaire		٠	٠	۰	۰	٠	٠		
4) Systemic Sclerosis Enrollment and	ACR Criteria - Lupus	۲	۲	۲	۲	۲	۲	۲		
Diagnosis Date Info 5) ACR	ACR Criteria- Scleroderma	۲	٠	٠	٠	٠	۰	٠		
6) females with older brothers	ACR Criteria - JIA	۲	۲	٠	۰	۲	۰	۲		
7) PLE SLEDAI	ECLAM	۲		٠	٠	٠	٠	۲		
 Plasma Cytokines-MMc JRA Rx 	SLEDAI		٠	٠	٠	٠	٠	٠		
10) MMc SLEDAI ECLAM	Lab Values	٠	٠	۲	٠	٠	٠	٠		
11) SN Cytokines ACR Renal	Genotyping	۲								

Fig. 2A.) REDCap is a research electronic data capture application used to track patient visit information and lab test results.

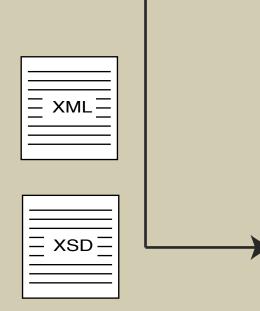


Fig. 2C.) Biomediator data integration system retrieves and caches X and schemas from rem applications.

schema, and answer

their patient data.

	Account	Listin	g que	eries				
	Welcome: demo		📥 XML	🖄 csv	Name	View	Edit	Delete
•	My Account Logout		\$7>		Full join of REDCap and rLab data	0		×
	Usage		\$7>		Aliquots for female patients w/older brothers needing DYS14	Q		×
a	Queries		\$7>		Aliquots for patients needing maternal qPCR and/or MLR	Q		×
ata	Actions		\$7>		Aliquots for patients needing any qPCR	Q		×
lla	List Queries		\$75		Aliquots for patients with incomplete genotyping	Q		×
		6	\$75		20 million+ PBMC for PLE and NOP	Q		×

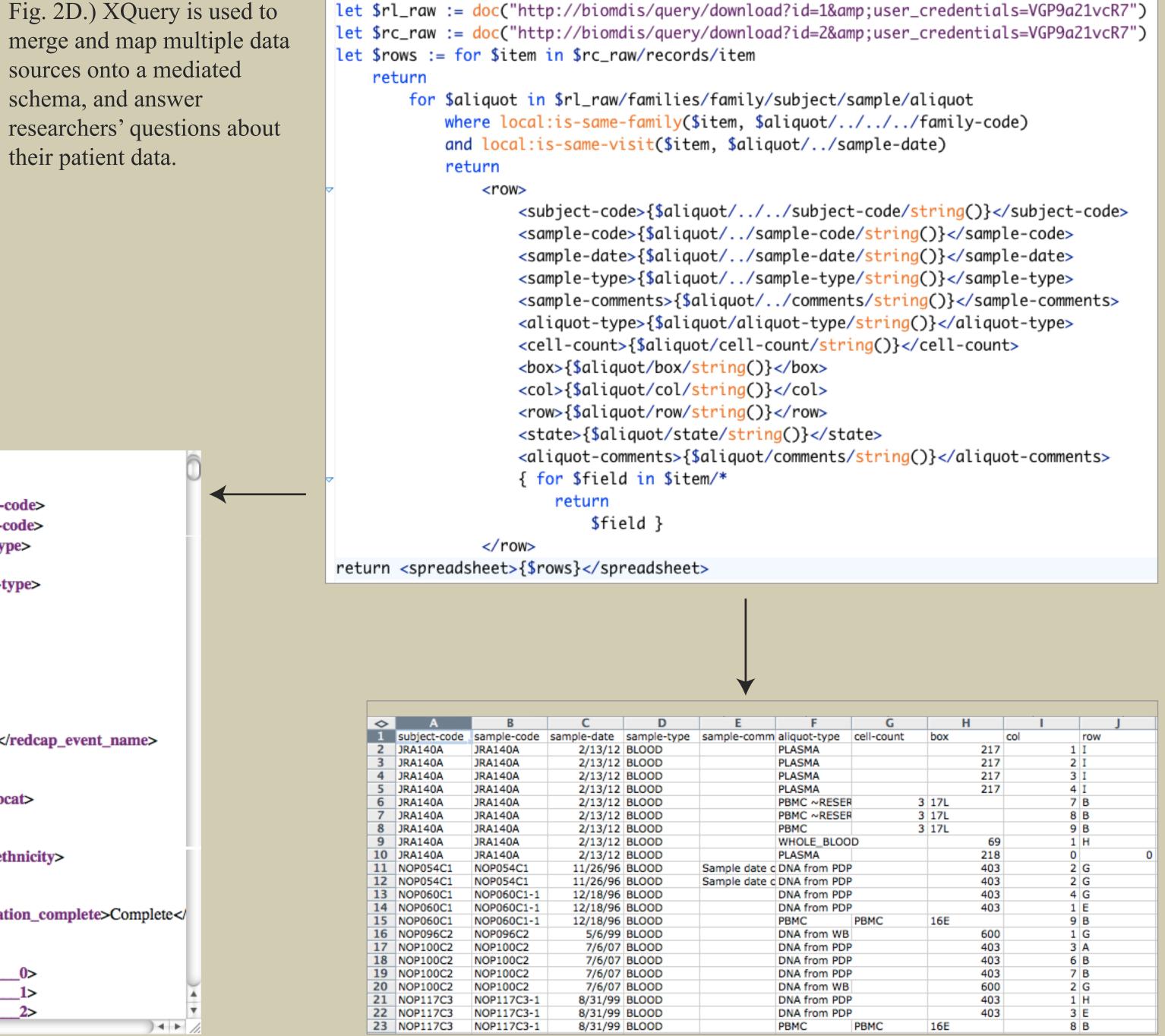


Fig. 2E.) XML result sets may be used to create a new data source that can be consumed by Biomediator and other data integration systems.

Fig. 2F.) Result sets may be transformed into other formats such as CSV (spreadsheet) and HTML.

8/31/99 BLOOD 8/31/99 BLOOD

22 NOP117C3 NOP117C3-1

NRM134 LIST VIEW SEARCH	H MY CLIPBOARD	SEARCH RESULTS LAB NEWS
EW FAMILY NEW SUBJECT EXPOR	RT (XLS)	QUICK SEARCH
Family: NRM134	Aliquot Info	
 Subject: NRM134 Sample: NRM134-7 	Aliquot type	PLASMA
	Aliquot label	2
	Box	49
	Row	D
	Col	2
-PLASMA: 49:D2	Cell Count	PLASMA millions
	State	NOT_USED
PBMC: 18D:E5	Aliquot used by	
	Aliquot used dat	e 0000-00-00
	Comments	

Fig. 2B.) rLab is a freezer manage application used to track patient sample types, availability and location.

	View	Edit	Delete		
	Q		×		
′S14	Q		×		
	Q		×		
	Q		×		
	Q		×		
	Q		×		
			cR7") cR7")		
ode))				
ole- ole- ole- mple	bjec code date type	> > > ment			



Use Case

To illustrate the use of BDIS, we've developed a small set of queries in support of a pediatrics auto-immunity study (Fig. 2C). The study's researcher stores their patient visit and lab data in an EDC application called REDCap (Fig 2A). The locations of patient blood samples taken during visits are stored in a separate freezer inventory system called rLab (Fig 2B). In order to answer questions such as "Are any samples available for patients who still need qPCR testing completed?", lab technicians were required to query both systems separately and merge those results by hand.

By translating their question into an XQuery across data downloaded from REDCap and rLab, manual reconciliation is avoided. In fact, by providing the end-user with a generic query and transformation that combines all data from both systems into a single spreadsheet (Fig 2D), we've found that end-users adept with programs like Microsoft Excel are able to answer previously difficult questions with relative ease.

Conclusion

The Biomediator Data Integration System enables researchers to more easily answer cross-application queries about their data. Our framework's inexpensive, minimalist approach to data integration provides re-usable set of building blocks for merging data from collections of heterogeneous sources.

Acknowledgements

Funded by NIH grants RR0254, AR051545, and NS055088.

References

Mork P, Halevy A, Tarczy-Hornoch P. A Model for Data Integration Systems of Biomedical Data Applied to Online Genetic Databases. Proceedings of the AMIA Annual Symposium; 2001 Nov 3-7; Washington D.C..