

A Survey of Informatics Platforms That Enable Distributed Comparative Effectiveness Research Using Multi-institutional Heterogenous Clinical Data

Dean F. Sittig, PhD,* Brian L. Hazlehurst, PhD,† Jeffrey Brown, PhD,‡ Shawn Murphy, MD, PhD,§ Marc Rosenman, MD,|| Peter Tarczy-Hornoch, MD,¶ and Adam B. Wilcox, PhD#

Abstract: Comparative effectiveness research (CER) has the potential to transform the current health care delivery system by identifying the most effective medical and surgical treatments, diagnostic tests, disease prevention methods, and ways to deliver care for specific clinical conditions. To be successful, such research requires the identification, capture, aggregation, integration, and analysis of disparate data sources held by different institutions with diverse representations of the relevant clinical events. In an effort to address these diverse demands, there have been multiple new designs and implementations of informatics platforms that provide access to electronic clinical data and the governance infrastructure required for interinstitutional CER. The goal of this manuscript is to help investigators understand why these informatics platforms are required and to compare and contrast 6 large-scale, recently funded, CER-focused informatics platform development efforts. We utilized an 8-dimension, sociotechnical model of health information

technology to help guide our work. We identified 6 generic steps that are necessary in any distributed, multi-institutional CER project: data identification, extraction, modeling, aggregation, analysis, and dissemination. We expect that over the next several years these projects will provide answers to many important, and heretofore unanswerable, clinical research questions.

Key Words: methods, comparative effectiveness research, organization and administration, medical informatics

(*Med Care* 2012;50: S49–S59)

The American Recovery and Reinvestment Act of 2009 provided \$1.1 billion for comparative effectiveness research (CER).¹ The goal of CER is to generate new evidence on the potential effectiveness, benefits, and harms of different treatments, diagnostics, preventions, and care models under “real world” conditions. Widespread adoption of CER has potential to radically change health care. CER also places enormous demands on existing informatics research infrastructure,² as it requires aggregation and analysis of disparate data held by different institutions, each with its own representation of relevant events and accountabilities for protecting data as a matter of patient confidentiality and business operations.

Currently, most data manipulations are performed using noncoordinated applications [eg, data collection forms, electronic health records (EHRs), research databases, condition-specific registries, and statistical analyses] with disjointed institutional control. In an effort to address these demands, there have been new designs and implementations of informatics platforms that provide access to electronic clinical data and the governance required for interinstitutional CER.^{3–6} Briefly, a “platform” is a suite of interconnected, coordinated applications, together with the operational environment that hosts those applications.

The goal of this manuscript is to compare and contrast 6 large-scale projects that are either developing or extending existing informatics platforms for CER. Rather than compare the informatics platforms at an abstract level, we focus on specific CER projects that provide implementations of informatics platforms and highlight design requirements and solutions.

The following sections provide an overview of the projects surveyed.

From the *School of Biomedical Informatics, University of Texas Health Science Center, Houston, TX; †Kaiser Permanente Center for Health Research, Portland, OR; ‡Department of Population Medicine, Harvard Pilgrim Health Care Institute; §Partners Healthcare System, Harvard Medical School, Boston, MA; ||Regenstrief Institute and Indiana University School of Medicine, Indianapolis, IN; ¶Division of Biomedical and Health Informatics, University of Washington, Seattle, WA; and #Department of Biomedical Informatics, Columbia University, NY.

D.F.S. is supported in part by a grant from the National Library of Medicine R01-LM006942 and by a SHARP contract from the Office of the National Coordinator for Health Information Technology (ONC #10510592); B.L.H.’s work is supported in part by grants from the National Library of Medicine (R21LM009728), and the Agency for Healthcare Research and Quality (AHRQ) (R01HS019828, R18HS18157); J.B. is supported in part by the AHRQ grant 1R01HS019912; S.M. is supported by grants from the National Institute of Health U54LM008748, UL1RR025758, U24RR025736 and by a grant from the Office of the National Coordinator 90TR0001/01; M.R. is supported by a grant from AHRQ (R01HS019818); P.T.-H. is supported in part by AHRQ 1 R01 HS 20025-01 “Surgical Care and Outcomes Assessment Program (SCOAP) Comparative Effectiveness Research Translation Network (CERTN)” and by NIH NCRR 1 UL1 RR 025014 “Institute of Translational Health Sciences”; A.B.W. is supported in part by AHRQ grant R01 HS019853-01, Washington Heights/Inwood Informatics Infrastructure for Community-Centered Comparative Effectiveness Research (WICER).

The authors declare no conflict of interest.

Reprints: Dean F. Sittig, PhD, UT-Memorial Herman Center for Healthcare Quality & Safety, 6410 Fannin St., UTPB 1100.43, Houston, TX 77030. E-mail: dean.f.sittig@uth.tmc.edu.

Copyright © 2012 by Lippincott Williams & Wilkins
ISSN: 0025-7079/12/5007-0S49

WASHINGTON HEIGHTS/INWOOD INFORMATICS INFRASTRUCTURE FOR COMPARATIVE EFFECTIVENESS RESEARCH (WICER)

WICER is creating infrastructure to facilitate patient-centered outcomes research in Washington Heights, NY. The project facilitates comprehensive understanding of populations by leveraging data from existing EHRs, and combining data from institutions representing various health care processes. For example, it includes data from hospitals, clinics, specialists, homecare agencies, and long-term care facilities. It also includes survey data from community residents with assessments on socioeconomic status, vital statistics, support networks, health and illness perceptions, quality of life, and health literacy. Data from multiple sources are merged in a data warehouse, where deeper analysis is performed by clinical and public health researchers. WICER investigators are using the infrastructure and methods on 3 clinical trials in hypertension care around diagnosis, adherence to therapy, and care management.

SCALABLE PARTNERING NETWORK FOR COMPARATIVE EFFECTIVENESS RESEARCH: ACROSS LIFESPAN, CONDITIONS, AND SETTINGS (SPAN)

The HMO Research Network (HMORN) is a consortium of 19 Health Plans with formal, research capabilities.⁷ SPAN, a project within the HMORN, uses its Virtual Data Warehouse (VDW) to provide a standardized, federated data system across 11 partners, to address CER in ADHD and obesity.⁸ The VDW consists of commonly defined linked tables within each health plan that capture medical care utilization, clinical data, health plan enrollment information, demographics, detailed inpatient and outpatient encounter information, outpatient pharmacy dispensing data, laboratory test results, and vital signs.⁹ The VDW is augmented with state and local cancer registry information on date and cause of death for health plan members. Each plan maintains control of individual VDW data files and does not have access to files held by other HMORN sites. All HMORN participants must be capable of running, without modification, SAS programs distributed by other sites to execute against their local VDW. SPAN is pioneering use of a new platform, PopMedNet, that facilitates creation, operation, and governance of multisite, distributed health data networks.¹⁰

ENHANCING CLINICAL EFFECTIVENESS RESEARCH WITH NATURAL LANGUAGE PROCESSING OF EHR DATA—CER-HUB

The CER-HUB is an Internet-based platform for conducting CER. A central function of CER-HUB is facilitating (through online, interactive tools) development of a shared, data processor library that can be downloaded by registered researchers to provide uniform, standardized coding of both free-text and structured clinical data. This shared library permits researchers to assess data on clinical effectiveness in multiple health care areas and gain access to information

locked in free-text notes. Using CER-HUB, researchers collaboratively build software applications (MediClass applications¹¹) that will process EHR data within their respective health care organizations, creating standardized datasets that can be pooled to address-specific CER protocols. Participating researchers contribute Institutional Review Board-approved, limited datasets to a centralized coordinating center to be pooled with data similarly processed from other health care organizations to answer CER questions. The CER-HUB is being used to conduct 2 CER studies addressing effectiveness of medication for controlling asthma and of smoking cessation counseling services, across 6 geographically distributed and demographically diverse health systems. Researchers and data providers for these initial studies come from 3 Kaiser health plans (Northwest, Hawaii, and Georgia regions), 1 consortium of Federally Qualified Health Centers located primarily along the west coast (OCHIN Inc.), 1 Veterans Administration service region (Puget Sound VA in Washington), and an integrated network of hospitals and physicians in the greater Dallas/Fort Worth area (Baylor Health Care System).

THE PARTNERS RESEARCH PATIENT DATA REGISTRY (RPDR)

The RPDR is an enterprise data warehouse combined with a multifaceted user interface that enables clinical research and CER across Partners Healthcare in Boston, MA. The RPDR is used to recruit patients for clinical trials, and to perform active surveillance. It amasses data from billing, decision support, and EHRs in the Partners' system. Data were available to researchers through a drag-and-drop web Query Tool¹² allowing users to construct exploratory, ad hoc, queries for hypothesis generation from structured data, and to get aggregate totals and graphs of age, race, sex, and vitals. A utility exists for finding matched controls for patients. Requests can be made for detailed data on patients identified through the query tool with proper IRB authorization through an automated wizard. The RPDR has proven useful for gathering clinical trial cohorts, and for CER. This strategy was later adopted as the core of "Informatics for Integrating Biology and the Bedside" (i2b2).¹³ The RPDR was first released in December 1999 and has been in production at multiple sites since March 2002.

THE INDIANA NETWORK FOR PATIENT CARE (INPC) COMPARATIVE EFFECTIVENESS RESEARCH TRIAL OF ALZHEIMER DISEASE DRUGS (COMET-AD)

INPC was begun in 1994 as an experiment in community-wide health information exchange (HIE) serving 5 major hospitals in Indianapolis, IN. Today, it includes data from hospitals and payers statewide.^{14–16} Entities participating in INPC submit patient registration records, laboratory test results, diagnoses, procedure codes, and other data for various types of health care encounters. Data are also obtained from health departments and a pharmacy benefit manager consortium. Data are standardized (eg, laboratory test results are mapped to LOINC¹⁷ with common units of

measure) to the extent possible, before storage in a central repository. Data for a patient with visits to multiple INPC institutions can be linked using a patient matching algorithm. The COMET-AD project is using data from INPC to monitor health care processes and outcomes and to build systems to monitor patients for adverse drug events. The project also involves building infrastructure and workflows to support integration of biospecimen results with clinical data from the INPC.

THE SURGICAL CARE OUTCOMES ASSESSMENT PROGRAM COMPARATIVE EFFECTIVENESS RESEARCH TRANSLATION NETWORK (SCOAP-CERTN)

The goal is to assess how well an existing statewide quality assurance and quality improvement registry (ie, the Surgical Care Outcomes Assessment Program) can be leveraged to perform CER. The SCOAP-CERTN leverages relationships built collaboratively in SCOAP to improve surgical care and outcomes and aims to build infrastructure for streamlined, electronic data abstraction from EHRs, patient-reported outcomes, and health care payments across hospitals. Through a partnership with Microsoft Health Solutions Group (Redmond, WA), SCOAP-CERTN is identifying ways to maximize automatic capture of data from EHRs, to:

- Allow longitudinal clinical data capture across health care encounter types (ie, surgical, interventional).
- Reduce clinical workflow and staffing burdens for maintenance of the SCOAP registry at participating hospitals.
- Provide capacity and interoperability to incorporate outpatient care delivery into SCOAP.

In addition, SCOAP developers plan to add functions to capture patient-reported outcomes for research and quality improvement evaluation. The primary informatics goal is to assess how, and to what degree, the collection of SCOAP-CERTN measures can be automated across sites.

CONCEPTUAL MODEL FOR CER PLATFORM EVALUATION

Designing, developing, implementing, and using health information technology within health care delivery systems is a complex, sociotechnical challenge. To provide a theoretical basis for our comparison of 6 CER informatics platforms we adapted an 8-dimension, sociotechnical model of safe and effective health information technology use.¹⁸ This model prescribes attention to: (1) appropriate hardware/software, (2) a spectrum of clinical content ranging from case narrative, to standard vocabularies, to algorithms representing best practices, (3) human-computer interfaces enabling productive interactions with technology, (4) personnel who develop systems and how systems meet the needs of users in their social contexts, (5) workflow and communications (both between people and technology components) required to accomplish tasks using the technology, (6) organizational policies, procedures, culture, and environment that prescribe and govern how and where things happen and who is responsible, (7) external rules, regulations, and

pressures which shape these organizational constraints, and (8) system measurement and monitoring which ensures adequate performance for primary intended use cases, that is, the conduct of CER.

These 8 constructs¹⁸ are used to investigate and evaluate aspects of CER platform design and implementation by ensuring that both the social as well as the technical aspects are considered. Failure to consider who will use the applications, how they will use them, and why they are necessary often leads to suboptimal technology design and utilization.

METHODS

Data Sources

We developed a written survey and sent it to informatics experts representing 6 large CER projects focusing on the design, development, and use of multi-institutional informatics platforms. Projects were selected by convenience, yet they are representative of vastly different approaches researchers have taken to address numerous CER challenges.

Survey Instrument

We (D.F.S., B.L.H.) developed a 2-page, open-ended survey that highlighted project-specific similarities and differences. We created 2–8 questions within each of the 8 dimensions to ensure that all important aspects were captured.¹⁸ For example, within workflow/communication we asked, “How do data get into your warehouse?” and “What stages do the data go through?” Similarly, within the hardware/software dimension we asked, “What computing infrastructure is required to run your system?”

Data Collection and Analysis

Completed surveys were returned by e-mail and checked for completeness. D.F.S. and B.L.H. read through the 6–10-page responses from each of the coauthors looking for key concepts highlighting project similarities and differences. After review and discussion, it became clear that the following 4 dimensions of the 8-dimension model were the key differentiators: content or data (Table 2); workflow/communication regarding how data moved from sources to analysis (Table 3); people (investigators, data programmers, research analysts, managers) involved in the projects (Table 4); and organizational policies, procedures, and culture (Table 5). We extracted data items to fill-in the tables from surveys. In addition to survey items, 2 authors (D.F.S., B.L.H.) gathered information regarding project descriptions and funding from websites and journal articles (Table 1). Drafts of completed tables were sent to coauthors for review.

RESULTS

All projects implement 6 generic data processing steps necessary for distributed, multi-institutional CER projects:

- Identification of applicable data within health care transaction systems.
- Extraction to a local data warehouse for staging.
- Modeling data to enable common representations across multiple health systems.

TABLE 1. Overview of the 6 Informatics Platform-based Comparative Effectiveness Research Projects Analyzed in this Manuscript

CER Project Name	Project PI/ Representative	Project-specific Funding ¹⁹	Geographic Diversity	Project Website	Areas of Study (No. Patients Involved)	No. Participating Institutions	Stage of the Infrastructure and Project
SPAN	Daley/Brown	\$8,272,272	USA-based Hawaii-Georgia	popmednet.org (software platform); Hmoresearchnet work.org has information on the VDW	Attention-deficit hyperactivity disorder (~240,000 adolescents); obesity (~2.4 million)	11 US-based health systems	Infrastructure: VDW in production 10y, SPAN is adding new variables; developing query enhancements to the PopMedNet platform Project: Developing research protocols, initial data assessments
WICER	Wilcox/Wilcox	\$8,855,607	Limited to Washington Heights/Inwood community (New York City, NY)	Wicer.org	Population in research data warehouse (~200,000); population cohort analysis using community surveys (12,000); care management, new diagnostic algorithms, 3 clinical trials of (600 patients total)	4 local organizations	Infrastructure: under development Project: in planning stages
CER-HUB	Hazlehurst/Hazlehurst	\$8,735,759	USA-based Hawaii-Georgia	cerhub.org	Asthma (200,000); smoking cessation (500,000)	6 large health systems—3 HMOs, 1 VA, 1 consortium of FQHCs, 1 large regional physician/hospital network	Infrastructure: under development Project: 2 studies underway
RPDR	Murphy/Murphy	~\$18 millions over past 12y	Limited to Boston, MA area	rc.partners.org/rpdr,12b2.org	Supports about 400 studies/year; 4 studies have been completed ^{20,21,25,26}	7 local hospitals	Infrastructure: System in operation since 1999 Project: Hundreds of studies underway
INPC COMET-AD	Dexter and Boustani/Rosenman	\$8,422,410	Indianapolis, IN metropolitan area	www.regenstrief.org/medinformatics/inpc	Drug treatment of Alzheimer's disease (300)	5 major hospital systems in Indianapolis area	Infrastructure: HIE used routinely for patient care and research Project: Beginning patient enrollment
SCOAP-CERTN	Flum/Tarczy-Hornoch	\$11,690,974	Washington state	www.scoop.org	Primary: peripheral artery disease; Secondary: under development	10 institutions planned (3 agreed)	Infrastructure: WA statewide QA/QI registry in use at 50+ sites for 4+ years Project: Recruiting 7 additional sites (3 agreed)

FQHCs indicates Federally Qualified Health Care facilities; HMO, Health Maintenance Organization; QA/QI, quality assurance/quality improvement; VDW, virtual data warehouse.

TABLE 2. Comparison of Data Sources, Types, Models, and Handling of Duplicate Patients

CER Project Name	Data Sources	Data Types	Standard Data Model(s) Used	Duplicate Patients Identified Across Organizations?
SPAN	Health plan enrollment, pharmacy dispensing, utilization data, billing data, vitals, lab results, tumor registry, death info	Local codes Standard codes No unstructured text	Expanded version of the HMO Research Network Virtual Data Warehouse Version 3 (13 tables linked by a unique identifier) ²²	No. Two organizations unlikely to have information for the same patient during a defined enrollment period.
WICER	Patient surveys, vital statistics, health literacy, socioeconomic status, in-patient, ambulatory clinics, long-term care facilities, home care agencies	Local codes Standard codes Processed free-text	Early version of the HL-7 Reference Information Model	Yes, many patients are participants in New York Care Connect HIE
CER-HUB	Ambulatory EHR, inpatient discharge summaries, billing, pharmacy dispensing, lab results; all are extracted based on project need by standard extraction mechanism	Local codes Standard codes Processed free-text	Implementation of HL-7 Clinical Document Architecture that extends the CCD (Continuity of Care Document)	No. Unlikely for sites currently involved to have overlap in patient populations. One site operates a single instance EHR for its multiple consortium member FQHC organizations.
RPDR	Demographics and labs data loaded nightly; EHR, billing and decision support systems data (including vitals and inpatient and ambulatory clinics data), death info and pharmacy data loaded monthly; text clinical notes available on project-need basis	Local codes Standard codes Processed free-text	Star schema data model codes clinical events as “facts” in relational database structure with radiating tables that further define facts, along with metadata tables	Yes, Enterprise Master Patient Index
INPC COMET-AD	Multiple hospital systems, health care payers, practice organizations (eg, primary care group practice, radiology practice, sports medicine practice), laboratory organizations	Local codes Standard codes Processed free-text and unstructured text	Identifier, timestamp, “question” term, and “answer” term—where answer term is numeric, coded, date, person (eg, patient or clinician), or free-text value. Also some “compound” results.	Yes, patients are linked across institutions in the Indiana Network for patient care through the global person ID service
SCOAP-CERTN	ADT/registration, laboratory, medications, text reports (eg, Doppler ultrasound report), text notes (eg, operative note)	Local codes Standard codes Processed free-text Unstructured text	HL7-v3 in warehouse augmented by data elements from the SCOAP data collection forms	No

FQHCs indicates Federally Qualified Health Care facilities; HIE, Health Information Exchange; HL7-v3, Health Level 7-Version 3.

- Aggregation of data according to this common data model.
- Analysis of data to address research questions.
- Dissemination of study results.

All projects performed these activities, although there were variations in how (real-time aggregation of HL-7 transactions vs. nightly or as-needed extraction, transformation, and loading), where (local site vs. coordinating center), and with what tools (web-based query interfaces for researchers vs. tools to develop natural language processing modules).

Table 2 compares data sources, types, models, and handling of duplicate patients. All projects collected data from multiple sources (ie, hospitals, clinics, billing, long-term care) and included different data types (eg, numeric test results, ICD-9-encoded problems, and free-text progress notes). Only 3 projects used a “master patient index” that enabled them to combine data from patients who received treatment at different organizations. All projects used different, and sometimes multiple, data storage and manipulation formats ranging from SAS tables to XML-based documents to relational databases.

Table 3 provides a comparison of data flow and transformation, from local EHRs to aggregated analyses. The most important differences highlighted in Table 3 pertain to when patient-identifiable data leave local sites. In 2 projects, this occurs immediately after extraction from the local transaction-based clinical or administrative systems. In SPAN and CER-HUB, transfer of “raw” patient-identifiable data never occurs (ie, all data were processed at the local site by data analysis programs that are distributed from the central site, and only data conforming to protocol-specific Limited Data Sets are shared). Only 3 sites had any form of natural language processing capability; the other sites relied solely on numeric or coded data elements.

Also of interest in Table 3 is the state of data analysis tools offered by projects. All projects are working on “user-friendly” tools to facilitate researchers’ direct access to data by ad hoc queries, whereas concurrently meeting multi-institutional requirements for protecting patient data and corporate business interests. To date, only the RPDR has a working version.

TABLE 3. Comparison of Data Flow and Transformation From EHR to Aggregated Analysis

CER Project Name	Raw Data	Natural Language Processing	Data Normalization	Data Aggregation	Data Analysis Tools
SPAN	Data from local EHR, billing, and other sources are accessed in their native form, or extracted and stored in research data warehouse	No	Each site transforms local data to the common data model (HMORN VDW) using native coding systems (NDC, ICD9, HCPCS, LOINC); transformations checked centrally for consistency	Centrally developed application (analytic programs, queries) runs at local site, results sent to Data Coordinating Center or study lead site	Menu-driven query tools for encounter-level VDW datasets are under development; tools available querying of aggregate data
WICER	Extracted from local EHRs, sent to central data warehouse where it is combined, by patient, with data from multiple organizations (HIE warehouse-model)	Yes, using general purpose MedLee ²³	Central DB uses Columbia's Medical Entities Dictionary (MED) ²⁴ to map to LOINC, ICD, SnoMed	Central site	Tool under development to allow end-users to identify patients with characteristics of interest, specify query constraints, and combine data elements
CER-HUB	Data extracted from Local EHR, stored in local data warehouse	Yes, using project-specific, knowledge-based MediClass applications ¹¹	Local site uses centrally developed data processor—MediClass application—to populate Clinical Research Documents using Unified Medical Language System as standard knowledge base	Centrally developed processors are downloaded to create study-specific clinical events that are transmitted to central site for aggregation and analysis	Tools available for researchers to create new and review existing NLP query modules and test them against new datasets. Tools in development to enable investigators to analyze and review final aggregated results.
RPDR	Data extracted from local systems (some daily, some monthly), aggregated centrally in data warehouse	Yes, using project-specific pattern recognition models ^{25,26}	Data mapped to ICD-9-CM and COSTAR (for diseases), NDC (for medications), LOINC (for labs), CPT and HCPCS (for procedures) data stored centrally in local coding systems and dynamically mapped in query tool	All aggregate data are derived from queries generated through query tool accessing central data repository, detailed data on patients are gathered from central repository and through enterprise web services	Drag-and-drop web Query Tool allows users to construct ad hoc, Boolean queries for hypothesis generation from structured data, to get aggregate totals and to graph age, race, sex, and vitals
INPC COMET-AD	Extracted from local EHR (or payer), sent to central data warehouse, stored distinctly but can be combined at patient-level across multiple organizations; HIE-model	Yes (for a few use cases thus far), using project-specific pattern recognition models. ²⁷	Data are normalized and coded using standard vocabularies (LOINC, ICD9, CPT, NDC, etc.) when/where appropriate	Specially developed data aggregation routines developed and run centrally	All SAS and other statistical data analysis done at the central site; new data analysis tools are under development
SCOAP-CERTN	Data are extracted from source systems and transmitted by VPN channels to the central data warehouse servers; message stream is queued then parsed	Text mining tools operate on the data within the central data warehouse; use hybrid text mining (rules base and statistical) to extract and tag the text and derive data elements needed for QA/QI and CER	Parsed data stream is either directly mapped onto the HL-7 data model or into a staging table from which it is mapped to the HL-7 derived and augmented data model	Data extracted from local systems and sent to secure remote server cluster running Amalga	Sites and SCOAP-CERTN personnel will use secure web sites and interfaces with appropriate authentication, authorization and logging to query the Amalga data warehouse

Gray cells indicate when data first leaves control of local site.

Table 4 describes key personnel. The most important difference is that some projects either have or are working on Internet-based interfaces that allow nontechnical investigators to perform a limited set of data queries and analyses

on the combined dataset. For example, the SPAN project currently requires all queries be coded as SAS programs and sent to the local site where they are executed and the results returned after manual review; SPAN is beta-testing an

TABLE 4. Key Personnel Involved in Various Stages of Project

CER Project Name	Data Extraction, Transformation, Loading	Data Queries	Data Analysis	Total Personnel Accessing System
SPAN	Local site programmers create HMORN VDW (multipurpose research warehouse)/expand VDW on study-by-study basis; distributed queries create study-specific analytic file extracts from the VDW	SAS analysts at the lead study site create SAS programs to run against the VDW. Nontechnical researchers use the SPAN user interface to query summary count tables; user interface enhancements will allow menu-driven querying of individual-level data derived from the VDW.	Statisticians, investigators, analysts at lead study site	Planned 20–40 investigators (2–4/site); 25 analytic/administrative staff to respond to queries (~2 per site)
WICER	Local site programmers work with central programmers to create extracts—real-time and monthly	Designed to be completed by researchers	Statisticians and investigators at local or central sites	Platform still under development
CER-HUB	Local site programmers create XML-based clinical research document	Study staff create and validate standardized data processors—a MediClass application—for each study	Statisticians at central data coordinating center with guidance from study investigators	Per study (2 currently) Local site programmers: 1–3 per study site (currently 12 total). Research staff: 5–10 investigators, RA's, and PM's per (currently 15 total) >2500 across the Partners' organization
RPDR	Local site programmers create real-time and periodic data extracts	Clinical investigators run queries using structured user interface	Preliminary statistics automatically generated by query tool. Investigators at study sites perform more detailed analysis.	
INPC COMET-AD	Central staff extracts data from the INPC; additional NLP is planned; data are also collected by project staff	SQL queries on central DB are written by central staff in collaboration with the investigator team	Analysis is by statisticians at the central location using extracted data	For the pilot study: 1 project manager and staff, 1 INPC data analyst, 1 NLP expert programmer, 6 investigators. The more general platform is under development
SCOAP-CERTN	Central staff work with local site programmers to create site-specific data extracts that are sent to centrally developed Amalga message processing architecture	Users at participating sites, centralized SCOAP-CERTN staff (site associated data coordinators, centralized data coordinators and scientists doing data analysis)	By statisticians at central location using extracted data	System is still under development. When ready, expect users across the state of Washington

NLP indicates natural language processing; SQL, Structure Query Language; VDW, Virtual Data Warehouse.

Internet-based approach using the PopMedNet architecture to allow nontechnical users to issue queries.

Table 5 provides a comparison of project governance and internal organizational policies and procedures. All projects have an oversight committee; most consisting of representatives from all sites involved in the project. Often this committee is responsible for governing all aspects of data ownership and sharing, project membership, and publication rights and responsibilities.

DISCUSSION

We compared 6 large CER projects and described how they use informatics platforms to provide data aggregation, analysis, and research management capabilities. Many of

these platforms were originally designed and developed to address widely different health care, organizational, and research objectives; only after significant amounts of work had been completed were they transitioned to focus on CER. For example, the RPDR was originally designed to answer the question, “How many patients with a specific set of characteristics have we treated within our integrated delivery network?” On the other hand, INPC and WICER started as a means of improving the quality and efficiency of care in large metropolitan areas by creating centrally managed HIEs. Similarly, SCOAP-CERTN started as a registry to improve surgical outcomes and efficiency. SPAN (and to a lesser extent CER-HUB) build upon existing research networks comprised of similarly organized and managed, large, integrated health plans.

TABLE 5. Comparison of Project Governance and Internal Organizational Policies and Procedures

CER Project Name	Project Governance	Project Membership	Data Ownership and Access Control	Data Use Agreements	IRB Oversight	Study Design Oversight	Publication Rights
SPAN	Project Governance Core, SPAN PI and executive team ²⁸	Inclusion in the project grant submission; additional membership possible by governance committee	Local sites control access to all data; opt-in on a study-by-study basis	Between local sites and coordinating center and study lead site(s); BAA agreement between sites and IT vendor that hosts the SPAN networking infrastructure	Each project overseen by local IRBs; most sites have ceded IRB oversight to the study lead site	All projects result from external contracts or grants. Each site must agree to participate.	Maintained by the study team in accordance with grant and contract language
WICER	Project Governance committee	All members were part of New York Care Connect HIE; new members by governance committee (community surveys are a new group, that were not represented by NYCC HIE)	PHI removed from data after matching patient IDs; all researchers have access to data through login ID	Local sites have DUA with central site. Each data source has a steward.	Local sites and central coordinating center all have IRBs	Researchers control study design; access to data through remote UI	Data steward defines right of access and first publication
CER-HUB	Steering committee (site PIs) authorize each project (currently 2 projects) but projects are autonomous teams addressing approved study protocol ²⁹	On the basis of projects run by a lead investigator; approved by Steering Committee	Raw data: stays at sites Study data: project-specific Limited Data Sets shared (using DUA) with data coordination center Results: project team controls with sites represented by an investigator	Between each site and the central data coordinating center	Sites each require an IRB submission to address local concerns and circumstance. One of the 5 secondary sites has ceded IRB oversight to the lead site.	Approved by steering committee through acceptance of project. Managed by study lead investigator.	Maintained by study team in accord with contract language and publications policy
RPDR	Project governance committee	On the basis of lead investigator always being a faculty member at the organization	Raw data: local sites Study data: stored centrally, controlled by query tool Results: local sites control, requires login managed by organizations involved	Established along with initial data transfer to central DB	Central IRB, enforced by query tool and requirement of digital signatures to obtain data extracts	Researchers control study design, user interface limits search strategy, and patient detail that may be viewed without explicit IRB approval	No oversight
INPC COMET-AD	INPC management committee (and/or its research subcommittee) oversees all data use; PI and coinvestigators lead the project; data and safety monitoring board for specific trials	All INPC members may participate (optional)	Raw data: local site owns data, but RI is data custodian Study data: stored centrally Results: managed by RI/Indiana University (or Purdue University or other) investigators and data analysts	Established between the INPC and each local hospital, or payer, etc	All research projects have separate IRB approvals, most with the Indiana University IRB	RI/university overseen by the University with initial data use approval from the INPC management/research committee/subcommittee	Controlled by the RI/Indiana University (or Purdue University or other) project investigator

SCOAP-CERTN	Overall External Advisory Boards including community representation (SCOAP and SCOAP-CERTN), also Health Information Technology-specific external advisory board	SCOAP-CERTN executive leadership identifies potential partners who are invited to participate	Sites retain control over their data; can elect to stop participation and withdraw their data from the central repository up until extraction for CER occurs	DUAs and BAAs being developed with IRB, compliance, Regulatory and Legal offices of the University of Washington and participating sites	Each substudy has an IRB protocol	Each CER project will oversee work	SCOAP-CERTN project publication guidelines overseen by SCOAP-CERTN executive committee
-------------	--	---	--	--	-----------------------------------	------------------------------------	--

BAA indicates Business Associate Agreements; DUA, Data Use Agreements; HMORN, Health Maintenance Organization Research Network; INPC, Indiana Network for Patient Care; RI, Regenstrief Institute.

CER Requires Comprehensive Data on Patients

Different data types are required to create complete, patient-centered views of patients' medical histories. The surveyed projects demonstrate that creating a useful CER platform requires enormous amounts, and a large variety, of data. To access these data, CER investigators need to collect them from as many different sources within their participating organizations as possible. Therefore, we see researchers collecting data from inpatient and outpatient EHRs (including the text narrative of clinical encounters), from billing and ancillary systems such as laboratory, pharmacy, and radiology. In addition, it is important to collect data that document that patients actually received the care that was ordered, so we see organizations collecting pharmacy dispensing and patient-reported data when available. This vast array of data, although large, is nearly always incomplete (ie, they generate sparse representations in a large-dimensional space of patient care facts in the real world) and methods which use these data must be appropriate to the task of measuring health status and care events with available data.

CER Requires Data on Populations From Multiple Organizations

Researchers need to aggregate data from multiple organizations to have enough information to identify small differences, address bias, perform subgroup analyses, improve generalizability, allow evaluation of demographic and geographic variation, and identify rare events. Therefore, CER informatics platforms must be able to extract and collect data from many different organizations to compile a complete view of conditions, treatments, and individuals as possible. Toward that end we see investigators working to include data from multiple organizations, pursuing non-traditional research data sources, such as long-term care facilities, home and public health agencies, and attempting to reliably ascertain patients' socioeconomic status on a widespread basis.

A key requirement for data collection across health care provider organizations located in the same geographic region is the need to merge data from the same patient who has received health care services and had clinical data captured at multiple institutions. Such efforts require a community-wide master patient index that identifies patients on the basis of multiple demographic data (eg, first name, last name, date of birth, sex, social security, or telephone numbers) and keeps track of all patient identifiers used by various participating organizations to create a single, master patient identifier.³⁰ To date, only the CER projects that were built on top of existing HIE platforms designed for patient care have tackled this extraordinarily difficult problem,^{31,32} but in the future patient matching capabilities will be a critical success factor.

CER Requires Data Extraction, Modeling, Aggregation and Analysis Methods and Tools

Researchers must be able to extract required data from various electronic data systems, map data types to standardized clinical representations, and analyze it. Design and development of these "mapping" applications is one of the

biggest challenges in any multi-institutional research project, because it is often the case that different organizations refer to the same activity, condition, or even procedure by different names, and the same names can refer to different things across institutions. Further, even with accurate mapping it is difficult for researchers to fully appreciate local idiosyncratic data issues (eg, nonrandom incomplete data capture) without active engagement of local data experts.

Furthermore, conducting CER is a complex undertaking requiring people with widely different skills, often in different locations and subject to different organizational policies and practices. In an attempt to reduce potential for misunderstanding in collaboration processes, platform developers are working to create powerful, user-friendly tools for data extraction, manipulation, and analysis. These tools are being designed so CER project staffs, who often have little informatics training, can perform their tasks more efficiently. In addition, several projects are developing tools to help researchers make sense of highly variable and clinically rich free-text notes documenting patient care.

CER Must Conform to Local Organizations' Internal Governance and IRBs' Rules and Local and Federal Legislation

The social, legal, ethical, and political challenges involved in setting up and conducting large, multi-institutional CER projects must not be underestimated. Friedman et al³³ stated that "organizations are understandably reluctant to move data beyond their own boundaries absent a clear and specific need to do so, and patients will be less likely to consent to allow this to happen." Therefore, in addition to providing the technical infrastructure required to collect, standardize, normalize, and analyze disparate data, informatics platforms must conform to local organizations' internal governance and IRBs' rules and regulations as well as existing state and federal guidelines. One design to address use of protected health information is to retain physical control of raw data although providing for their aggregation as limited datasets to answer specific questions. Other ways in which projects have accommodated interinstitutional governance issues include standardizing data models across the project; limiting access to authorized personnel while facilitating remote access; restricting the types of queries that can be executed and masking patient-specific, identifiable data; and logging all data transactions and access activities. As rules, regulations, and guidelines evolve (eg, proposed Common Rule revision³⁴) CER platforms and governance processes must evolve accordingly.

SUMMARY AND CONCLUSIONS

CER stands to transform the current health care delivery system by identifying which therapies, procedures, preventive tests, and health care processes are most effective from the standpoints of cost, quality, and safety. State-of-the-art informatics platforms are necessary to carry out this type of research across organizations with disparate patient populations, health information systems, data types, and local governance structures.

We used an 8-dimension, sociotechnical model to develop a survey enabling us to compare and contrast informatics platforms that are under development or in use in 6 large CER efforts. On the basis of the data we collected, we identified 6 generic steps necessary in any distributed, multi-institutional CER project: data identification, extraction, modeling, aggregation, analysis, and dissemination.

We conclude that all of the informatics platforms for CER studied are on their way to creating the sociotechnical infrastructure required to enable researchers from multiple institutions to conduct high-quality, cost-effective CER. We expect that over the next several years, these projects will provide answers to many important CER questions that in the past were virtually inaccessible. In addition, we expect many more CER-focused informatics research platforms to be designed, developed, and tested as the fields of informatics and CER continue to evolve.

ACKNOWLEDGMENT

The authors thank Andrea Bradford, PhD, for editorial assistance.

REFERENCES

1. The American Recovery and Reinvestment Act of 2009. Public Law 111-5. Available at: <http://www.gpo.gov/fdsys/pkg/PLAW-111publ5/pdf/PLAW-111publ5.pdf>. Accessed February 17, 2009.
2. VanLare JM, Conway PH, Rowe JW. Building academic health centers' capacity to shape and respond to comparative effectiveness research policy. *Acad Med*. 2011;86:689–694.
3. Stang PE, Ryan PB, Racoosin JA, et al. Advancing the science for active surveillance: rationale and design for the Observational Medical Outcomes Partnership. *Ann Intern Med*. 2010;153:600–606.
4. Ohno-Machado L, Bafna V, Boxwala AA, et al, and the iDASH team iDASH: integrating data for analysis, anonymization, and sharing. *J Am Med Inform Assoc*. 2012;19:196–201.
5. Behrman RE, Benner JS, Brown JS, et al. Developing the Sentinel System—a national resource for evidence development. *N Engl J Med*. 2011;364:498–499.
6. Payne P, Ervin D, Dhaval R, et al. TRIAD: The Translational Research Informatics and Data management grid. *Appl Clin Inf*. 2011;2:331–344.
7. Greene SM, Hart G, Wagner EH. Measuring and improving performance in multicenter research consortia. *J Natl Cancer Inst Monogr*. 2005;35:26–32.
8. Toh S, Platt R, Steiner JF, et al. Comparative-effectiveness research in distributed health data networks. *Clin Pharmacol Ther*. 2011;90:883–887.
9. Hornbrook MC, Hart G, Ellis JL, et al. Building a virtual cancer research organization. *J Natl Cancer Inst Monogr*. 2005;35:12–25.
10. Brown JS, Holmes JH, Shah K, et al. Distributed health data networks: a practical and preferred approach to multi-institutional evaluations of comparative effectiveness, safety, and quality of care. *Med Care*. 2010;48(suppl 1):S45–S51.
11. Hazlehurst B, Frost HR, Sittig DF, et al. MediClass: a system for detecting and classifying encounter-based clinical events in any electronic medical record. *J Am Med Inform Assoc*. 2005;12:517–529.
12. Murphy SN, Gainer VS, Chueh H. A visual interface designed for novice users to find research patient cohorts in a large biomedical database. *AMIA Annu Symp Proc*. 2003:489–493. Available at: www.ncbi.nlm.nih.gov/pmc/articles/PMC1480150/pdf/amia2003_0489.pdf. Accessed May 4, 2012.
13. Murphy SN, Weber G, Mendis M, et al. Serving the enterprise and beyond with informatics for integrating biology and the bedside (i2b2). *J Am Med Inform Assoc*. 2010;17:124–130.
14. Overhage JM, Tierney WM, McDonald CJ. Design and implementation of the Indianapolis Network for Patient Care and Research. *Bull Med Libr Assoc*. 1995;83:48–56.
15. McDonald CJ, Overhage JM, Barnes M, et al. INPC Management Committee. The Indiana network for patient care: a working local health

- information infrastructure. An example of a working infrastructure collaboration that links data from five health systems and hundreds of millions of entries. *Health Aff (Millwood)*. 2005;24:1214–1220.
16. Zhu VJ, Tu W, Rosenman MB, et al. Facilitating clinical research through the health information exchange: lipid control as an example. *AMIA Annu Symp Proc*. 2010;2010:947–951.
 17. McDonald CJ, Huff SM, Suico JG, et al. LOINC, a universal standard for identifying laboratory observations: a 5-year update. *Clin Chem*. 2003;49:624–633.
 18. Sittig DF, Singh H. A new socio-technical model for studying health information technology in complex adaptive healthcare systems. *Qual Saf Healthc*. 2010;19(suppl 3):i68–i74.
 19. Health & Human Services Research Awards: Use of Recovery Act and Patient Protection and Affordable Care Act Funds for Comparative Effectiveness Research. U.S. Government Accountability Office, Washington, D.C.; June 14, 2011. Available at: <http://www.gao.gov/new.items/d11712r.pdf>. Accessed May 4, 2012.
 20. Tatonetti NP, Denny JC, Murphy SN, et al. Detecting drug interactions from adverse-event reports: interaction between paroxetine and pravastatin increases blood glucose levels. *Clin Pharmacol Ther*. 2011;90:133–42. doi: 10.1038/clpt.2011.83. [Epub 2011 May 25].
 21. Kurreeman F, Liao K, Chibnik L, et al. Genetic basis of autoantibody positive and negative rheumatoid arthritis risk in a multi-ethnic cohort derived from electronic health records. *Am J Hum Genet*. 2011;88:57–69.
 22. CHAPTER 4: VIRTUAL DATA WAREHOUSE (VDW) In Collaboration Toolkit: A guide to multicenter research in the HMO Research Network. 2011. Available at: www.hmoresearchnetwork.org/resources/toolkit/HMORN_CollaborationToolkit.pdf#4.
 23. Friedman C, Shagina L, Lussier Y, et al. Automated encoding of clinical documents based on natural language processing. *J Am Med Inform Assoc*. 2004;11:392–402.
 24. Baorto D, Li L, Cimino JJ. Practical experience with the maintenance and auditing of a large medical ontology. *J Biomed Inform*. 2009;42:494–503.
 25. Zeng QT, Goryachev S, Weiss S, et al. Extracting principal diagnosis, comorbidity and smoking status for asthma research: evaluation of a natural language processing system. *BMC Med Inform Decis Mak*. 2006;6:30.
 26. Perlis RH, Iosifescu DV, Castro VM, et al. Using electronic medical records to enable large-scale studies in psychiatry: treatment resistant depression as a model. *Psychol Med*. 2012;42:41–50. [Epub 20 June 2011].
 27. Friedlin J, McDonald CJ. Using a natural language processing system to extract and code family history data from admission reports. *AMIA Annu Symp Proc*. 2006:925.
 28. HMO Research Network Collaboration Toolkit. Available at: http://www.hmoresearchnetwork.org/resources/collab_toolkit.htm.
 29. CERHUB Mission and Membership Overview. Available at: <http://www.cerhub.org/mission.html>.
 30. AHIMA, Fundamentals for Building a Master Patient Index/Enterprise Master Patient Index (Updated), *J AHIMA*. Available at: http://library.ahima.org/xpedio/groups/public/documents/ahima/bok1_048389.hcsp?dDocName=bok1_048389. Accessed September 2010.
 31. Adragna L. Implementing the enterprise master patient index. *J AHIMA*. 1998;69:46–48, 50, 52.
 32. McDonald CJ, Overhage JM, Tierney WM, et al. The Regenstrief medical record system: a quarter century experience. *Int J Med Inform*. 1999;54:225–253.
 33. Friedman CP, Wong AK, Blumenthal D. Achieving a nationwide learning health system. *Sci Transl Med*. 2010;2:57cm29.
 34. Regulatory changes in advanced notice of proposed rulemaking: comparison of existing rules with some of the changes being considered. Available at: <http://www.hhs.gov/ohrp/humansubjects/anprmchangeable.html>. Accessed September 13, 2011.