

Building an integrated neurodegenerative disease database at an academic health center

Sharon X. Xie^{a,b,c,d,e,*}, Young Baek^{b,d,e,f}, Murray Grossman^{c,d,e,f}, Steven E. Arnold^{d,e,g,h},
Jason Karlawish^{b,d,e,h,i}, Andrew Siderow^{d,e,f,j,k}, Howard Hurtig^{d,e,f,j,k}, Lauren Elman^{d,e,f,l},
Leo McCluskey^{d,e,f,l}, Viviana Van Deerlin^{b,d,e,j,m}, Virginia M.-Y. Lee^{d,e,m},
John Q. Trojanowski^{b,d,e,j,m}

^aDepartment of Biostatistics and Epidemiology, University of Pennsylvania School of Medicine, Philadelphia, PA, USA

^bAlzheimer's Disease Core Center, University of Pennsylvania School of Medicine, Philadelphia, PA, USA

^cCenter for Frontotemporal Degeneration, University of Pennsylvania School of Medicine, Philadelphia, PA, USA

^dInstitute on Aging, University of Pennsylvania School of Medicine, Philadelphia, PA, USA

^eCenter for Neurodegenerative Disease Research, University of Pennsylvania School of Medicine, Philadelphia, PA, USA

^fDepartment of Neurology, University of Pennsylvania School of Medicine, Philadelphia, PA, USA

^gDepartment of Psychiatry, University of Pennsylvania School of Medicine, Philadelphia, PA, USA

^hPenn Memory Center, University of Pennsylvania School of Medicine, Philadelphia, PA, USA

ⁱDivision of Geriatrics, Department of Medicine, University of Pennsylvania School of Medicine, Philadelphia, PA, USA

^jMorris K. Udall Parkinson's Disease Research Center of Excellence, University of Pennsylvania School of Medicine, Philadelphia, PA, USA

^kParkinson's Disease and Movement Disorder Clinic, University of Pennsylvania School of Medicine, Philadelphia, PA, USA

^lAmyotrophic Lateral Sclerosis Center, University of Pennsylvania School of Medicine, Philadelphia, PA, USA

^mDepartment of Pathology and Laboratory Medicine, University of Pennsylvania School of Medicine, Philadelphia, PA, USA

Abstract

Background: It is becoming increasingly important to study common and distinct etiologies, clinical and pathological features, and mechanisms related to neurodegenerative diseases such as Alzheimer's disease, Parkinson's disease, amyotrophic lateral sclerosis, and frontotemporal lobar degeneration. These comparative studies rely on powerful database tools to quickly generate data sets that match diverse and complementary criteria set by them.

Methods: In this article, we present a novel integrated neurodegenerative disease (INDD) database, which was developed at the University of Pennsylvania (Penn) with the help of a consortium of Penn investigators. Because the work of these investigators are based on Alzheimer's disease, Parkinson's disease, amyotrophic lateral sclerosis, and frontotemporal lobar degeneration, it allowed us to achieve the goal of developing an INDD database for these major neurodegenerative disorders. We used the Microsoft SQL server as a platform, with built-in "backwards" functionality to provide Access as a frontend client to interface with the database. We used PHP Hypertext Preprocessor to create the "frontend" web interface and then used a master lookup table to integrate individual neurodegenerative disease databases. We also present methods of data entry, database security, database backups, and database audit trails for this INDD database.

Results: Using the INDD database, we compared the results of a biomarker study with those using an alternative approach by querying individual databases separately.

Conclusions: We have demonstrated that the Penn INDD database has the ability to query multiple database tables from a single console with high accuracy and reliability. The INDD database provides a powerful tool for generating data sets in comparative studies on several neurodegenerative diseases.

© 2011 The Alzheimer's Association. All rights reserved.

Keywords:

Database; Neurodegenerative disease; Microsoft SQL; Relational neurodegenerative disease database

Sharon X. Xie and Young Baek are co-first authors.

*Corresponding author. Tel.: 215-573-3867; Fax: 215-573-4865.

E-mail address: sxie@mail.med.upenn.edu

1. Introduction

Age-related neurodegenerative diseases, such as Alzheimer's disease (AD), Parkinson's disease (PD), amyotrophic lateral sclerosis (ALS), and frontotemporal lobar degeneration (FTLD), are increasing in prevalence with the rise in longevity of populations in most countries across the globe. For example, members of the “baby boomer” generation in the United States (i.e., all those born between the years of 1946 and 1964) turned 60 in January 2006, and, in 2011, will turn 65 [1]. Notably, AD increases exponentially after the age of 65 with its prevalence doubling every 5 years [1]. However, increasing evidence shows that AD and most other neurodegenerative diseases (including PD, FTLD, and ALS) share similar types of hallmark pathologies (e.g., disease-specific protein aggregates) and there exists an overlap in the clinical features of these disorders. For example, Neumann et al [2] showed that TDP-43 is the disease-causing protein in both ALS and FTLD; in addition, it is well-known that TDP-43 pathology occurs in PD and AD, as well as in other disorders to a variable extent [3,4]. Therefore, it is critical that we conduct multidisciplinary patient-oriented clinical and basic science research comparatively and in a multidimensional manner to improve the understanding and treatment options for diseases such as AD, PD, FTD, ALS, and other aging-related neurodegenerative disorders. Moreover, advances in one of these disorders could accelerate the pace of advances for other such diseases. To achieve this goal, it is essential to build an overarching integrated neurodegenerative disease (INDD) database, which includes multiple neurodegenerative disorders such as AD, PD, FTD, ALS, and other related diseases. Specifically, using this neurodegenerative disease database as a research tool, investigators would be able to obtain data across several disease groups and conduct comparative studies to elucidate distinct and common features and mechanisms of these disorders.

In this article, we describe a novel neurodegenerative disease database, which we developed at the University of Pennsylvania (Penn) with the help of a consortium of Penn investigators working collaboratively on AD, PD, FTLD, and ALS. In the Methods section, we describe the technical details of building this relational neurodegenerative disease database. In the Results section, we present snapshots of the database, provide detailed information regarding the types of data captured by the database, and using the INDD database, compare the results of a biomarker study with those using an alternative approach by querying individual database separately. In the Discussion section, we consider the advantages, limitations, challenges, and benefits of constructing the INDD database. We have also provided the development costs and the anticipated cost savings. Finally, we discuss how this technology will aid in new discoveries, our future development plans, and the potential effect on neurodegenerative disease research.

2. Methods

First, a flowchart indicating the steps and timelines of building the INDD database is provided in Table 1. These steps have been elaborated in the following subsections.

2.1. Step 1: Preparation and meeting to discuss criteria of INDD database

From the inception of database implementation, in which paper forms used for data capture were converted to electronic databases, each individual neurodegenerative disease research group at Penn which had unique focuses on AD, PD, FTLD, or ALS used locally-housed databases. These databases comprised various database platforms and technologies, and ranged from group accessible Microsoft Access (Version 2007; Microsoft Corporation, Redmond, WA) [5] databases and single-user Microsoft Excel (Version 2007; Microsoft Corporation) spreadsheets to more complex SPSS files and FileMaker databases (SPSS Inc, Chicago, IL). In 2006, the leaders of these different disease-centered research groups at Penn embarked on a joint, consensus-driven effort to create a single Penn INDD database. The most mature database at the time was developed through the National Institute of Aging funded Alzheimer's Disease Core Center (ADCC), and the first step in creating a single comprehensive database was to determine common technologies and methods that could be used for converting the individual databases into a new integrated database.

The first stage of conceptualizing the integrated database required two initial criteria that had to be addressed by each neurodegenerative disease research group or center before the launch of the INDD database project. The first criterion was that the new database had to be “backwards” compatible with the older database. More specifically, this meant that all the data from the new database had to have an ability to export to Excel files and then be accessible and queryable using the Access interface because some of the centers were accustomed to using an Access database.

The second criterion was that locally-assigned patient identification numbers (IDs) were to be not altered. Each individual center's patient IDs were linked and used on various

Table 1
Flow chart of steps in building the INDD database

Step 1 (2006) Preparation and meeting to discuss criteria of INDD database
Step 2 (2006) Selecting backend of the database to meet criterion 1: new database had to be “backwards” compatible with the old database
Step 3 (2006) Determining the programming language to create the frontend
Step 4 (2006) Integration of database to meet criterion 2: locally assigned IDs were not to be altered
Step 5 (2006) Building a rapid database development environment
Step 6 (2007) Implementing a comprehensive database security system
Step 7 (2007) Implementing database audit trails
Step 8 (2007) Implementing database backups
Step 9 (2007) Incorporating various data entry methods
Step 10 (2007) Implementing quality control procedures

Abbreviation: INDD, integrated neurodegenerative disease.

paper charts and patient samples, including biomarker samples, because requiring each center to renumber all patients and related items was deemed as an ineffective strategy.

2.2. Step 2: Selecting backend to meet criterion 1: new database had to be “backwards” compatible with the old database

When deciding on the platform upon which to develop and implement the new database, both of the aforementioned criteria had to be taken into consideration, and the collective final decision was to implement INDD with the Microsoft SQL (MSSQL) server (Version 2005; Microsoft Corporation) [6]. In 2006, when this project was launched, MSSQL was a widely deployed relational database management system with built-in backward functionality to provide Access as a frontend client to interface with the database. Although the primary user interface for the new integrated database was based on dynamic web-developed technologies, the MSSQL server provided an easy implementation of backwards compatibility with Access to satisfy the first criterion.

2.3. Step 3: Determining the programming language to create the frontend

After selecting the “back end” relational database management system component, the next step was to determine the programming language to create the “frontend” web interface, and we decided to develop the new database by using the PHP Hypertext Preprocessor (PHP; Version 5.2.3.) [7]. PHP is a widely deployed dynamic web language specifically created for developing web pages with flexible and powerful built-in functions that allow for quick access to a comprehensive online database. Along with PHP, other fundamental web technologies were used to develop the database, such as JavaScript [8], cascading style sheets [9], HTML [10], and XML [11]. Building on the web technology foundation, several other modern technologies have also been used, including asynchronous JavaScript and XML [12], which is a combination of several technologies, to create dynamic, user-friendly interfaces. In addition, jQuery [13] was used for rapid JavaScript development and an in-house built PHP framework geared toward the rapid development of patient-oriented database forms. The combination of these various technologies yielded powerful forms for accommodating various uses of the database.

2.4. Step 4: Integration of databases to meet criterion 2: locally assigned IDs were to be not altered

As stated previously, the INDD server comprised four clinical core or disease-centered databases along with three supporting databases (Fig. 1). The four clinical core databases (AD, PD, ALS, and FTLT) were created before the integration, thus the structures of the databases vary slightly between one disease-focused group or center and another by means of data collection. Although AD, PD, FTLT,

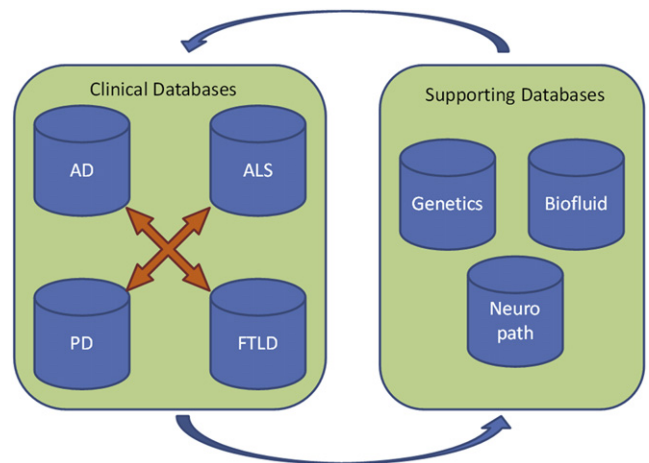


Fig. 1. Integrated neurodegenerative disease (INDD) database.

and ALS centers at Penn have their own unique center-specific IDs, it is entirely possible for a patient to be seen in multiple Penn centers. In this case, the patient will have several center-specific IDs. The three supporting databases are the biofluid database, neuropathology database, and genetics database. These databases comprised their own specific IDs with ties to the clinical center IDs. All seven databases are linked by the use of a master lookup table. Each unique individual added to any of the clinical databases will have a unique entry in this lookup table. When information on a new patient is requested in one of the clinical databases, the system will perform a lookup on the basis of the patient's name. The database will search for a similar-sounding name on the basis of the built-in MSSQL function SoundEx [14]. By using SoundEx, similar sounding names are used to display already existing patients in the database. Once identified, the database displays a list of similar patients and prompts the user to choose an already-existing patient with the verification of the date of birth and address to form a patient linkage through a unique entry in the master lookup table. This allows different database users to be aware of the changes in the information of a patient on the other databases. A special warning system is in place to generate alerts to the different clinical centers when patient information is updated from a separate center.

Although the database maintains the information of individual patients in an accurate state through a lookup table and global database IDs, when data are requested and queried from the database, overlapping IDs from different centers pose difficulties in identifying the individuals with ease and accuracy. To address this problem, the database automatically appends a prefix letter to distinguish between originating centers in cases when a multicenter query is performed with a potential for overlapping patient IDs. For the four clinical centers, the individual IDs are identified with a composite ID constructed with a capital letter in addition to their patient ID. The method used in this instance is the capital letter A for AD center, P for PD center, F for FTLT center, and M for ALS center.

The last segment for completion of the integrated system was incorporating the genetics database into the INDD database. The genetics database is unique in that it uses its own separate third-party database system called Progeny (Version 7, Progeny Software, LLC, South Bend, IN), which has been designed specially for gathering and generating pedigrees from patient information housed in the database. Because of this advanced feature, we wanted to integrate our clinical database along with the Progeny database. Although Progeny is a closed source system, it allows access to the core database fields through open database connectivity. Using the provided open database connectivity, we are able to select, insert, update, and delete records to and from Progeny directly from our integrated database, thereby facilitating the integration of the genetics core with the clinical cores. This provided us with the necessary ability to include genetic information in queries across the database and to have Progeny “push” genetic information into the clinical database so that clinicians who have been cleared for access and have a validated credential can view patient genetic information. Because patient information and especially genetic information is highly confidential, a special genetics table was implemented with an extra security level. In addition to the standard researcher view group, a clinician view group was established, thereby allowing for only the primary clinician of the patient to be able to view the genetic data. This, along with a separation of the genetic data from the rest of the data in the INDD database, ensures that only a few investigators with appropriate clearance are privileged to view sensitive genetic information.

2.5. Step 5: Building a rapid database development environment

When developing an in-house customized database, especially a web-based database from a blank slate, it is important to be able to rapidly develop forms that can be used by the researchers. A hybrid research and clinical database, such as the INDD database, requires new forms to be created in conjunction with new studies. For the purposes of code reusability and rapid database development, an in-house database framework has been specifically created for the Penn INDD database. The framework is flexible enough to accommodate the diverse formats of the individual forms and also comprehensive enough for quick development, with a uniform look and feel across the database and forms. The INDD framework is based on model-view-controller architecture, where the framework includes MSSQL for data or as a model, PHP as controller, and HTML and cascading style sheets for presentation or view. There are various instances of usage of asynchronous JavaScript and XML for making direct database calls from JavaScript through PHP, thereby preventing unpleasant refreshing of the web pages while creating more of the look and feel of a desktop application.

2.6. Step 6: Implementing a comprehensive database security system

Because we were developing a medical research database containing sensitive patient information, security was a priority of utmost importance. To prevent access by unauthorized users, data transmission used encrypted tunnels, and several group- and user-level authorizations were implemented for data protection and authorization. The first layer of security, commonly called network layer security, is the initial line of defense against unauthorized activity. Network layer security is intended to position the database physically and logically in a secure location, thereby minimizing the exposure of the database to unauthorized users and machines. The INDD database is located in a University of Pennsylvania Health System (UPHS) secured network infrastructure and data center. The UPHS data center is located within the Hospital of the University of Pennsylvania and it houses all of the UPHS servers and data equipments. The data center is an environment-controlled, physically secured area that only allows authorized UPHS information technology personnel physical access to the servers. UPHS also provides a secure UPHS network, segregated from the Internet through a UPHS firewall and a secure network implementation. Access to the database is only possible while being physically connected to the UPHS network or through virtual private network [15] access to the UPHS network. Because the database server is located in the UPHS data center at Hospital of the University of Pennsylvania and protected by the UPHS firewall, the database has the same first tier level protection as the rest of the UPHS data (Fig. 2).

The second layer of security implementation, also known as domain layer security, involves users being authenticated on the operating system or the machine level of the database. Because the MSSQL is a Microsoft product and the UPHS data center is built on Microsoft domain architecture, one must have a user name on the UPHS domain to properly access the database server on the domain level.

The next level of security implementation is database layer security. Even with proper credentials to the server,

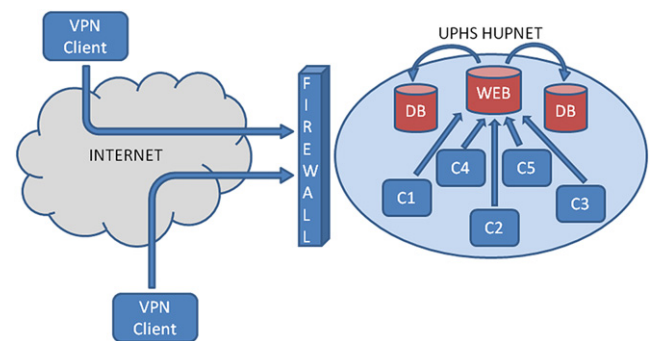


Fig. 2. System layout of the INDD database. (C1, client 1; C2, client 2; C3, client 3; C4, client 4; C5, client 5; DB, database; Web, web server; VPN client, virtual private network client; UPHS HUPNET, University of Pennsylvania Health System and Hospital of the University of Pennsylvania Network).

one must also obtain database level security approval to access the database. This layer provides admittance to the web-driven database along with globally defined table, and one can access the database by logging in using their database credential.

Once a user is authenticated to the main database container, each individual table and database form contain further security levels. Tables are configured with the following four data transaction operations: SELECT, INSERT, UPDATE, and DELETE. Each table is associated with one or more groups, thus restricting access to tables based on the group membership of a given user. Users are associated with groups, and the groups are applied to each individual table with the aforementioned four data transaction operations. The security groups are divided into individual centers, and then subdivided into the following four smaller subgroups: principal investigator, researcher, data entry, and read-only. The principal investigator security level allows SELECT, INSERT, and UPDATE operations on all center-specific tables, thereby allowing full control over the data fields. The researcher level allows SELECT, INSERT, and UPDATE on a subset of the tables. The data entry level allows access to even smaller sets of tables designed for capturing data transcribed from paper. Read-only access allows users to view the permitted information but without the ability to make any modifications. None of the aforementioned security groups are allowed to delete records from the database; a special database administration group is in place to allow data deletions. This was implemented to prevent accidental data deletions and to encourage good data entry habits. The decision regarding who has access to the database is made by a committee comprising the co-authors and collaborators in the author list who decide collectively on access on the basis of the qualifications of the investigators who request access and the nature of their project.

Data transmission between the server and the client is encrypted and secured using 128-bit transport layer security [15], also known as secure socket layer cryptographic protocols. Transport layer security prevents the occurrence of eavesdropping on the communication between the server and the client. This type of security is similar to the one implemented by banks and other institutions with sensitive data communications.

The INDD database uses a stringent data security model to thwart and prevent unauthorized data access and to minimize accidental data modifications. Along with the technical security model, security protocols, such as health insurance portability and accountability act and information technology policies, are in use for good practice and to encourage safe habits.

2.7. Step 7: Implementing database audit trails

An audit trail is one of the most important aspects of a reliable database. Irrespective of how well the database is designed with proper access levels, data constraints and process handling, the end result is ultimately still controlled by and dependent on human actions. The best defense against human errors is to ensure that proper backup and audit trails are in place. The INDD database relies on the third-party product ApexSQL Logs software (ApexSQL LLC, Chapel Hill, NC) for MSSQL server 2005. ApexSQL Logs provide comprehensive up-to-date audit trail information based on built-in MSSQL logs. ApexSQL Logs parse the log files to determine the information to be presented with numerous filters incorporated. Because ApexSQL Logs perform audit trails based on an already-in-place process of MSSQL, it has no overhead performance degradation or any requirement for additional storage. Figure 3 is an ApexSQL screenshot displaying activities in the INDD database for selected days.

State	Begin Time	Operation	Schema	Table	User	Client Host	Application	End Time
Committed	2009-08-07 12:07:25.143	Insert	dbo	PDC_Visits	sorvig	NULL	PHP 5	2009-08-07 12:07:25.143
Committed	2009-08-07 12:09:56.286	Insert	dbo	PDC_Visits	sorvig	NULL	PHP 5	2009-08-07 12:09:56.286
Committed	2009-08-05 15:57:34.003	Update	dbo	PDC_Visits	sorvig	NULL	PHP 5	2009-08-05 15:57:34.003
Committed	2009-07-28 15:45:01.750	Insert	dbo	PDC_Visits	sorvig	NULL	PHP 5	2009-07-28 15:45:01.750
Committed	2009-08-06 16:00:33.093	Insert	dbo	PDC_Visits	sorvig	NULL	PHP 5	2009-08-06 16:00:33.093
Committed	2009-08-05 15:57:13.300	Update	dbo	PDC_Visits	sorvig	NULL	PHP 5	2009-08-05 15:57:13.300
Committed	2009-08-05 15:52:10.543	Insert	dbo	PDC_Visits	sorvig	NULL	PHP 5	2009-08-05 15:52:10.543
Committed	2009-08-05 15:53:27.653	Insert	dbo	PDC_Visits	sorvig	NULL	PHP 5	2009-08-05 15:53:27.653
Committed	2009-08-05 15:54:44.090	Insert	dbo	PDC_Visits	sorvig	NULL	PHP 5	2009-08-05 15:54:44.090
Committed	2009-08-05 15:50:28.556	Insert	dbo	PDC_Visits	sorvig	NULL	PHP 5	2009-08-05 15:50:28.556
Committed	2009-08-26 12:09:01.950	Insert	dbo	PDC_Visits	sorvig	NULL	PHP 5	2009-08-26 12:09:01.950
Committed	2009-08-06 16:02:14.690	Insert	dbo	PDC_Visits	sorvig	NULL	PHP 5	2009-08-06 16:02:14.690
Committed	2009-08-06 14:45:25.546	Insert	dbo	PDC_Visits	sorvig	NULL	PHP 5	2009-08-06 14:45:25.546

Field	Type	Value	NewValue
PDCID	int	2298	2298
VisitDate	datetime	2009-07-28 00:00:00.000	2009-07-28 00:00:00.000
HNSE	float	30	30
HandYon	float	NULL	2
HandYoff	float	NULL	NULL
MedSinemet	bit	0	0
MedPernax	bit	NULL	NULL
MedAntidepressant	bit	0	0
MedAmantadine	bit	0	0
MedSinemetCR	bit	0	1
MedMirapex	bit	0	1
MedSeroquel	bit	0	0
MedClozapine	bit	0	0

Fig. 3. ApexSQL screenshot (database audit screenshot).

2.8. Step 8: Implementing database backups

Along with the audit trail, the INDD database is also backed up by using traditional database backups. The INDD database server is part of the UPHS enterprise backup system, and all information on the UPHS network, including the INDD database, undergo backup daily through an automated tape library system. Additionally, a separate database backup is performed on a weekly basis. This is primarily for information archiving and this backup is kept offsite in a permanent archive.

2.9. Step 9: Incorporating various data entry methods

In the INDD database, three different data entry methods are in use to accommodate different data and research needs. These data entries were performed at each local clinic. The first method is paper form data entry; in this entry method, a database user enters the data through an electronic data entry form that links directly to the database. This paper-capturing method is used when a computer is not directly available for the data to be captured or when the situation does not warrant the use of a direct computer entry. The electronic data entry form mimics the look and feel of the paper forms, thereby allowing for an easy transition from the paper to the database with minimal training and computer skill. The second data entry method allows for fast and accurate data entry by usage of dropdown numerical fields with number keypads along with tabs. On the basis of our experiences, this allows for a trained data entry person to perform speedy data entry with relatively high accuracy. The downside of this method is that personnel must be trained on how to perform the data entry. The third data entry method is importation of data. There are instances when the desired data are already stored in an Excel file or other delimited text files. When these files are to be imported to the database, we use an import function within the database. Data import allows existing data files to be directly transferred to the database, thereby bypassing user data entry. The benefit of this method is that data entry error is minimized because no human transcription is required. However, the downside of this method is that the file being imported must be predefined and correctly formatted, and an error could either halt the process altogether or incorrect data could be entered without the knowledge of the user.

Incorporating various data entry methods provides flexibility and usability to the database. Maximum efficiency can be achieved by selecting the data entry methods that are most suited to the situation at hand.

2.10. Step 10: Implementing quality control procedures

The Penn INDD database uses stringent data quality checks to ensure the accuracy and quality of the data. Double entry of randomly selected data enables us to monitor the quality of data entry using paper forms and we randomly se-

lect 10% of the original source records each quarter and perform a double data entry check. The double entry error rate is defined by the ratio of the total number of errors to the total number of data values doubly entered. Although our goal is 100% accuracy, an error rate of <0.5% is acceptable in general; anything higher will initiate a review of the data entry procedures (e.g., checking whether someone in particular has problems with entry or the data entry form is poorly designed) and appropriate corrective action is taken. Inconsistencies are reported immediately to the appropriate core/project leader for resolution.

Other quality control procedures include a wide range of data checking through range checks, consistency checks to ensure that data entered in the database are consistent with the expected data, as well as missing value and cross form consistency checks. We implement hard stops and soft stops in the data entry process so as to ensure these quality control procedures. A hard stop prevents a user from entering the data value in the database when the value does not fit within the range or the type expected. In contrast, a soft stop allows a user to enter the data value in the database and save it. However, the database will flag the value as questionable and require a user or an administrator to go back and double check on the data entered. At that time, either the data can be corrected by the user or the warning can be cleared by an administrator.

The database has routine scheduled jobs for gathering and providing reports of missing and/or questionable data. This report is regularly e-mailed out to the end users ensuring proper communication of data correction. Along with data quality control, the database also has the ability to communicate with a third-party reference for ensuring data accuracy. For example, the database communicates with the United States Postal Service address verification server to ensure that addresses in the database, including doctor offices and patient addresses, are correct. Along with internal procedures and processes to ensure data accuracy and quality, the database is flexible enough to have modules for communicating with external sources as well.

2.11. Validation of the INDD database

Using the Penn–Pfizer collaborative biomarker study as an example [16], we compared our INDD database approach to a traditional database approach with separate and disjointed database containers. We have compared the detailed steps involved in each method and presented the advantages using the INDD database versus the other approach in the Results section.

3. Results

As discussed in the Methods section, there are seven upstream databases within the INDD database: AD, PD, ALS, FTLT clinical databases, as well as biofluid database, neuropathology database, and genetics database. The INDD database includes measures in the areas of demographics,

clinical assessment, neuropsychological tests, imaging, pathology, biofluid, genetics, and clinical trials. Table 2 provides a summary of several variables in each of these arrays and some key variables. As of July 27, 2010, there were a total of 460,000 observations (unique records) in the INDD database.

Since the inception of the INDD database, there have been many examples of its usefulness and benefit in data retrieval, analysis, and research. This is most clearly illustrated by a recent series of biomarker-targeted proteomic studies that were performed across all disease domains in the INDD including AD, PD, FTLT, and ALS, as reviewed recently by Hu et al [16]. Because these studies included interrogation of approximately 1500 biofluid samples from several 100 patients using a multiplex system to measure >150 analytes in each sample, it is hard to imagine how we could have completed these studies without using the Penn INDD. Therefore, having a cross-disease database incorporating major neurodegenerative diseases (i.e., AD, PD, FTLT, ALS) along with biofluid samples, neuropathology, and genetic information has conferred great advantages in the quantity and quality of neurodegenerative disease data sets at Penn. As summarized in the review by Hu et al [16], abundant data fields within the database, as well as compatible data fields from across neurodegenerative disease centers, provided us with the information that was needed to correlate these biomarker data with clinical features of the different disorders. Therefore, these studies illustrate the exceptional data mining capabilities of the INDD database. Figure 4 provides an example of the INDD database interface with patient background and family history.

One of the best examples showcasing the advantages and strengths of the INDD database was a biomarker study conducted at Penn through the Penn–Pfizer Alliance in which 1500 plasma and cerebrospinal fluid (CSF) samples from patients with AD, PD, FTLT, or ALS and normal controls were interrogated using the Rules Based Medicine Inc., human discovery/multi-analyte panel of 151 analytes configured for the multiplex Luminex platform. The study initially required queries of the INDD database to ensure Penn had the necessary data from the four clinical disease centers to match various study criteria, as well as the ability to locate

and extract the corresponding plasma and CSF samples. The study criteria for selecting the cases required for a subject to have had either a plasma or CSF sample drawn from one of the four clinical centers with emphasis on having both plasma (e.g., epidermal growth factor) and CSF biomarkers (e.g., CSF *t*-tau). Additionally, each patient was required to have had a full clinical evaluation performed and had psychometrics tests (e.g., Mini Mental State Examination [MMSE]), vitals (e.g., blood pressure), and medical history (e.g., stroke) gathered.

We compared and contrasted two database methods to extract the data that satisfied the aforementioned criteria in the Penn–Pfizer collaborative biomarker study. We have demonstrated further how two different database schemes even after differing in design still manage to arrive at the same results.

The first database method used to generate the data was the traditional database design with separate and disjointed database containers. In this design, each clinical center housed their own center data locally using their center-specific IDs. Among others, a biofluid database, a neuropathology database, and a genetics database were also implemented in their individual containers segregated from others. When performing the same data extraction required by the aforementioned criteria for the Penn–Pfizer biomarker study, each of the four clinical center's databases were queried separately along with three supporting databases. Once the data were queried and the Excel data files were collected, the next step was to compare each of the files and ensure that no duplicate records of patients among different centers were found and then combine the four separate Excel files. In this post-processing of the data, one must carefully examine the data to ensure that no duplicate records are found and take extra care when combining the files. In this example, after querying the databases, the resulting data set contained more than 5000 records, which had to be examined and stitched together during post-processing. In a large study like this biomarker study, the investigators commonly request either the data to be rerun with additional data fields or to be rerun in the future after additional data have been added to the database. With this method of separate databases and the need to perform post-processing of the

Table 2
Summary of data fields and arrays of the INDD database

Array name	Number of variables	Key variables
Demographics	731	Education, race, gender
Clinical assessment	5830	Diagnosis, age of onset
Neuropsychological measures	3333	Mini Mental State Examination, Clinical Dementia Rating, Geriatric Depression Scale
Imaging	341	Volumetric data, imaging diagnosis
Pathology	646	Neuropathology diagnosis, amygdala tau, middle frontal gyrus tau
Biofluid	389	CSF <i>t</i> -tau, CSF <i>p</i> -tau, CSF Abeta
Genetics	117	<i>APOE</i> , <i>MAPT</i>
Clinical trials	8085	Consent date, visit date

Abbreviation: CSF, cerebrospinal fluid.

The screenshot displays the INDD database interface. At the top, there is a header for the University of Pennsylvania Neurodegenerative Diseases, with navigation tabs for MAIN, CLINICAL, RESEARCH STUDIES, NEUROPSYCH, UDALL, and REPORTS. Below the header is a search bar with fields for Patient ID, First Name, and Last Name, and a Search button. The main content area is divided into two sections: 'Background' and 'Family History'.

Background Section:

- Physician: [Dropdown]
- Site: [Dropdown]
- First Visit Date: [Date Picker]
- Year of Onset: [Text]
- Year of Diagnosis: [Text]
- DBS Date: [Date Picker]
- DBS Target: ☐ STN ☐ GPI ☐ VIM
- Symptom 1: [Dropdown]
- Symptom 2: [Dropdown]
- Symptom 3: [Dropdown]
- Other: [Text]
- ☐ Deceased
- Date of Death: [Date Picker]
- ☐ Autopsy
- Age at Death: [Text]

Family History Section:

- ☐ Patient has NO known family history of PD or Dementia.
- ☐ Patient has 1st Degree Relative with Dementia. Specify: [Dropdown]
- ☐ Patient has 1st Degree Relative with Parkinson's. Specify: [Dropdown]

Fig. 4. INDD database interface example.

data, the task of rerunning the data extraction is time consuming and challenging. All the steps of extracting and combining the data must be repeated for each instance, leaving room for human error and possible misrepresentation of the data.

The second method used to perform the data extraction was the INDD database method. Using the INDD database and its capability of centralized jointed tables, a single query was crafted to join 13 separate tables using the aforementioned criteria. The query generated 1103 records with each row representing a unique patient with the data points spanned across the columns. This result was exported to Excel, formatted and annotated for each column header, and then sent to the investigators for their analysis. In the event of rerunning the query, the INDD database stores previous executed queries in the database. Because the data extraction was performed through a single query, the query could be modified to contain the additional fields required by the investigators or the same query could be rerun to update the records of the data set.

In the previously mentioned case study, one can clearly see the advantage of the INDD database versus the traditionally deployed databases. The reduction in time and effort in using the INDD database enables researchers and data managers to focus their efforts elsewhere and eliminate the steps required by manual post-processing, thereby greatly reducing the chances of error in the data. Although the conclusions of the two data sets are identical, the two different approaches vary significantly in the time, effort, and accuracy of the resulting data set. Table 3 summarizes the key differences of the two database approaches.

With the ability to query across multicenter data sets and to match those data with biofluid and/or genetic data, the INDD database played a key role in our ability to conduct this study. Figure 5 gives an example of a portion of the data set queried for the Penn–Pfizer biomarker study from several clinical core databases. It shows data from the ALS, AD, FTD databases with education, race, ethnicity, and diagnosis along with MMSE date, MMSE total score, Luminex total CSF tau (t -tau) values, and Luminex CSF phosphorylated tau (p -tau) values. Since the completion of the interrogation of these 1500 plasma and CSF samples, several analyses of the data have been either already published, submitted, or are in preparation. Briefly, several analytical strategies are being used to identify important analytes on the basis of clinical and pathological diagnosis, including significance analysis of microarrays and random forest analysis. There were differences observed in many analytes obtained from subjects with AD and normal controls, whereas differences were seen in only a few analytes between those with AD and non-AD dementias. This type of analysis required for the model to adjust for basic demographic variables (age, gender, education) at the most superficial level, and additional adjustment for more complex time-dependent variables including duration of disease at the time of biofluid sample(s) collection and cognitive and neurological examination results corresponding to biofluid collection. Because some patients had multiple types of biofluids collected (plasma and CSF), and a small subset had serial samples collected at different time points, a comprehensive INDD database was required to generate the data points associated with each patient at a particular time point.

Table 3
Comparison between the INDD database approach and the separated database approach

Stages of generating data sets	INDD database approach	Separated database approach
Initial data planning stage	Step 1: Determine the necessary data values to be included in the query.	Step 1: Determine the necessary data points to be included in the query.
Data identifying stage	Step 2: With a database administrator, identify the tables containing the data points.	Step 2: With a database administrator, identify the databases and tables containing the data points.
Data gathering and formatting for analysis	Step 3: Write a single query that pulls from multiple tables interconnected by a key master table.	Step 3: Write a separate query from the individual database; a total of four queries have to be written; output the result to an Excel or other delimited file format. Step 4: Examine the individual files to ensure no duplicate records are found across different files. Step 5: Using Excel or other spreadsheet programs, join the separate files into a single large file. Careful attention is required to ensure the accuracy of matching up the data from individual files.
Rerunning of the query	Step 4: To rerun the query or run a modified query, adjust the single query then rerun.	Step 6: To rerun the queries or run modified queries, Steps 3 to 5 must be performed.

Novel analytes representing potential CSF biomarkers for AD and FTLD using the data generated from the INDD database have been studied previously and results have been either published [17] or submitted (Hu et al, unpublished data, 2010). We also investigated plasma biomarkers that distinguish subjects with AD from normal controls and other neurodegenerative diseases and these studies are being prepared for publication (Soares et al, unpublished data). Therefore, we exploited the Penn INDD database to implement novel biomarker studies that would otherwise have been fairly impossible to accomplish in a timely manner without an integrated database system.

4. Discussion

We developed an INDD database that included AD, PD, FTLD, ALS, and normal control patients, as well as those with other neurodegenerative diseases. This is significant because the INDD database provides a powerful tool for generating data sets in comparative studies across several neurodegenerative diseases.

In the Penn–Pfizer collaborative biomarker study, having an integrated database significantly eases the process of querying and extracting data from the database compared with the traditional database scheme, in which individual centers and individual components of the centers operate in separate and disjointed databases. We showed that the INDD database is robust and reliable because the two approaches generated identical results of the study. However, using a traditional database scheme is prone to a higher degree of errors and significant increase in labor for post-processing of the data. In contrast, the Penn INDD database has the ability to query multiple database tables from a single console with high accuracy and reliability. Therefore, the merits of this Penn INDD database are now evident and its usefulness as a research tool will certainly grow as it continues to mature and expand.

We kept and made separate distinctions between different centers and tables for the ease of data entries and not for technical reasons. The INDD database avoids altering local data entry. Although the interfaces for each center are customized and center-specific, all the data are stored in a same container and can be queried using joint databases.

	PatientID	Education	Race	Ethnicity	Diagnosis	TestDate	MMSETotal	Luminexttauave	Luminexptauave
1	M1018	15	White	Latino	ALS	02/13/2007	28	109.07	80.86
2	F17008	14	White	Non-Latino	CBD	03/18/2003	21	118.64	39.25
3	M1172	16	White	Non-Latino	ALS	03/12/2007	23	378.6	145.65
4	M1311	16	White	Non-Latino	ALS	11/01/2007	28	231.08	124.26
5	A11.33	12	White	Non-Latino	MCI - memory plus other	10/05/2006	26	24.5	8.75
6	F6429	18	White	Non-Latino	AD	04/20/2009	26	919.647	121.946
7	F6431	16	White	Non-Latino	AD	03/03/2009	28	1545.972	187.834
8	M1070	12	White	Non-Latino	ALS	01/04/2007	30	147.28	94.85
9	A11.33	12	White	Non-Latino	MCI - memory plus other	11/06/2006	24	24.5	8.75
10	M1054	12	White	Non-Latino	ALS	04/27/2007	29	164.05	42.79
11	F17004	18	White	Non-Latino	CBD	04/05/2003	5	187.59	64.59

Fig. 5. Screenshot database output for Penn–Pfizer biomarker study. This screenshot shows data queried from the INDD database on parameters for patients with ALS, AD, or FTLD including education, race, ethnicity, and diagnosis along with Mini-Mental State Examination (MMSE) date, MMSE total score, Luminex CSF *t*-tau values, and Luminex CSF *p*-tau values.

One challenge in building the INDD database is to make sure that the subject present in more than one center is indeed the same person. The work presented in this article illustrates the use of SoundEx to display similar sounding patients and the use of the date of birth and address to verify that they are indeed the same subject. This feature is important for an integrated database that links multiple clinical centers because it is not unusual to have a subject who is seen in more than one clinic. The second challenge is to continually use the separate center IDs that were created before the integration took place while accommodating for instances where the center IDs could overlap with other centers. After multiple meetings among data managers and investigators, we came up with the idea of building the INDD database by automatically appending a prefix letter to distinguish between originating centers. This is essential for researchers who plan to perform any comparative studies across more than one disease. The third challenge is to unify the variable coding. For example, race variable may have different categories for different centers. Data managers and principle investigators had to meet to decide on the common coding scheme for each common variable. This is obviously crucial because any comparative study and statistical analysis would require the same definition of the same variable for all subjects.

The current limitation of the INDD database is that the data queries have to be run by data managers and cannot be run by investigators who request the data. Our future development plan is to build a query bank that includes common queries that researchers can run by themselves.

In terms of the cost of the development, we purchased a new server with SQL software and license (about \$6000). It took approximately 75% effort of a data manager per year for about 2 years to build the INDD database. However, these investments will eventually save costs because all centers will share one server instead of using their own servers. Furthermore, investigators do not need to clean and merge data sets by themselves when they need data sets from multiple centers. They can get a clean data set from the INDD database. This will greatly reduce errors and save time. Finally, as exemplified by the biomarker studies mentioned previously and reviewed by Hu et al [16], the INDD will be enormously useful for future biomarker studies, as well as for genomic and genome-wide association studies in addition to the analysis of clinical trials and classic clinicopathological correlations. Indeed, we view the INDD as a significant step forward in creating the foundational framework on which we and others can usher in the exciting new era of “personalized medicine.”

Acknowledgments

There are no financial or other relationships that could be interpreted as a conflict of interest affecting this manuscript. This work was supported by NIH grants including the AG-10124 (the ADCC), AG-17586, AG-09215, NS-44266,

AG-15116, and NS-053488 (the Morris K. Udall Parkinson's Disease Research Center of Excellence) and an anonymous foundation. The authors are grateful to all their colleagues in the Center for Neurodegenerative Disease Research, the Penn Memory Center, the Center for Frontotemporal Degeneration, the Parkinson's Disease and Movement Disorder Clinic, the Morris K. Udall Parkinson's Disease Research Center of Excellence, the ADCC, and the Amyotrophic Lateral Sclerosis Center for their contribution in making the Penn INDD database a reality. The authors also thank Warren Bilker, Jonas Ellenberg, John Holmes, Kurt Brunden, and William Hu for their helpful comments on the draft of the manuscript. Finally, they thank all their patients and their families for their participation in their research and for their support for the Penn neurodegenerative disease research programs.

References

- [1] Trojanowski JQ, Arnold SE, Karlawish JH, Brunden K, Cary M, Davatzikos C, et al. Design of comprehensive Alzheimer's disease centers to address unmet national needs. *Alzheimers Dement* 2010; 6:150–5.
- [2] Neumann M, Sampathu DM, Kwong LK, Traux A, Micsenyi M, Chou TT, et al. Ubiquitinated TDP-43 in frontotemporal lobar degeneration and amyotrophic lateral sclerosis. *Science* 2006;314:130–3.
- [3] Geser F, Martinez-lage M, Kwong LK, Lee VM-Y, Trojanowski JQ. Amyotrophic lateral sclerosis, frontotemporal dementia and beyond. *J Neurol* 2009;256:1205–14.
- [4] Pesiridis GS, Lee VM-Y, Trojanowski JQ. Mutations in TDP-43 link glycine rich domain functions to amyotrophic lateral sclerosis. *Hum Mol Genet* 2009;18:R156–62.
- [5] Forte S, Howe T, Wall K. Access 2002 Development. Indianapolis, IN: Sam Publishing; 2001.
- [6] Rankins R, Bertucci P, Jensen P. Microsoft SQL Server 2000 Unleashed. Indianapolis, IN: Sams Publishing; 2003.
- [7] Schlossnagle G. Advanced PHP Programming. Indianapolis, IN: Sams Publishing; 2004.
- [8] Goodman D, Morrison M. JavaScript Bible. Indianapolis, IN: Wiley Publishing Inc; 2007.
- [9] Bartlett K. Teach Yourself CSS. Indianapolis, IN: Sam Publishing; 2001.
- [10] Castro E. HTML 4 for the World Wide Web. 4th ed. Berkeley, CA: Peachpit Press; 2000.
- [11] Myer T. No Nonsense XML Web Development With PHP. Collingwood, Victoria, Australia: Sitepoint; 2005.
- [12] Holzner S. Ajax Bible. Indianapolis, IN: Wiley Publishing, Inc; 2007.
- [13] York R. Beginning JavaScript and CSS Development With jQuery. Indianapolis, IN: Wiley Publishing, Inc; 2009.
- [14] Turley P, Wood D. Transact-SQL With SQL Server 2000 and 2005. Indianapolis, IN: Wiley Publishing, Inc; 2006.
- [15] Dean T. Network+ Guide to Networks. 3rd ed. Boston, MA: Course Technology; 2004.
- [16] Hu WT, Chen-Plotkin A, Arnold SE, Grossman M, Clark CM, Shaw LM, et al. Biomarker discovery for Alzheimer's disease, frontotemporal lobar degeneration, and Parkinson's disease. *Acta Neuropathol* 2010;120:385–99.
- [17] Hu WT, Chen-Plotkin A, Arnold SE, Grossman M, Clark CM, Shaw LM, et al. Novel CSF biomarkers for Alzheimer's disease and mild cognitive impairment. *Acta Neuropathol* 2010;119:669–78.