

Informationsfusion — Herausforderungen an die Datenbanktechnologie

— Kurzbeitrag —

Stefan Conrad¹ Gunter Saake² Kai-Uwe Sattler²

¹ Institut für Wirtschaftsinformatik, Johannes Kepler Universität Linz
Altenberger Str. 69, A-4040 Linz, Österreich
conrad@dke.uni-linz.ac.at

² Fakultät für Informatik, Otto-von-Guericke-Universität Magdeburg
Postfach 4120, 39016 Magdeburg
{saake|kus}@iti.cs.uni-magdeburg.de

Zusammenfassung In vielen Anwendungsbereichen besteht die Aufgabe, Daten oder Informationen aus verschiedenen, zum Teil heterogenen Quellen zu kombinieren, zu verdichten und daraus Informationen einer neuen Qualität abzuleiten. Wesentliche Kernfunktionen dieses als *Informationsfusion* bezeichneten Prozesses sind dabei durch Methoden der Datenintegration und der Datenanalyse / Data Mining bereitzustellen. Die gewachsenen Strukturen der heute genutzten Informationsquellen und die damit im Zusammenhang stehenden Probleme wie Heterogenität, Inkonsistenz oder Ungenauigkeit der Daten sind mit den aktuell verfügbaren Techniken nur bedingt beherrschbar. Ausgehend vom aktuellen Stand der Forschung diskutiert der vorliegende Beitrag Anforderungen an Datenbanktechnologien aus Sicht der Informationsfusion, zeigt mögliche Forschungsrichtungen auf und skizziert aktuelle und zukünftige Anwendungsfelder.

1 Motivation

Der heutige Stand der Datenbanktechnologie ermöglicht die effiziente Speicherung und Verwaltung von Datenbeständen mit unterschiedlichen Strukturen im Giga- und Terabyte-Bereich. Gleichzeitig erlauben moderne Kommunikationsmedien wie das Internet den Zugriff auf weltweit verteilte Informationen. Als eine Folge dieser Entwicklungen sehen sich jedoch viele Anwender einer wachsenden Informationsflut ausgesetzt, die das Auffinden von relevanten Informationen erschwert. Verschärfend wirkt dabei auch, daß viele Datenbestände zum Teil anarchisch gewachsen sind und demzufolge heterogen (sowohl bezüglich der Struktur als auch der Repräsentation) sind sowie Redundanzen und Inkonsistenzen enthalten. Die Integration externer Informationen (z.B. aus dem World Wide Web) eröffnet zwar neue Nutzungspotentiale für die unternehmensinternen Informationssysteme, ist aber gleichzeitig mit neuen Fragestellungen verbunden, wie z.B. der Gewährleistung von Aktualität und Vertrauenswürdigkeit von

Diese Arbeit wurde teilweise vom Land Sachsen-Anhalt unter FKZ 1987A/2527R und von der DFG unter Sa 465/9 gefördert.

Informationen, dem effiziente Zugriff auf weltweit verteilte Quellen oder der Aufbereitung von unzureichend strukturierten Daten.

Darüber hinaus enthalten die integrierten Datenbestände oft auch Informationen, die nicht explizit abgelegt sind, sondern sich in Form von Abhängigkeiten, Beziehungen oder Mustern über die einzelnen Quellen hinweg repräsentieren. Bei der Suche und Extraktion dieser impliziten oder „versteckten“ Informationen versagen jedoch klassische Anfragetechniken aus dem Datenbankbereich. Eine wichtige Anforderung an Informationssysteme ist demzufolge eine (semi-)automatisierte und intelligente Transformation der Daten in „nützliche“ Informationen. Der Begriff der Transformation umfaßt dabei verschiedene Aspekte, wie Integration, Filterung, Analyse und Aufbereitung der Daten mit dem Ziel des Aufdeckens und der Repräsentation des impliziten Wissens.

Vor diesem Hintergrund sind neben dem Datenmanagement und der Datenintegration die Bereiche Data Mining und Data Fusion Gegenstand aktueller Forschung. *Data Mining* beschäftigt sich als Kern des Prozesses der Wissensfindung in Datenbanken (Knowledge Discovery in Databases - KDD) mit der Suche nach Mustern und Abhängigkeiten in Daten. *Data Fusion* beschreibt die Kombination und Interpretation von Daten aus verschiedenen Quellen. Der Einsatz dieser Techniken im Rahmen von Informationssystemen eröffnet neue Möglichkeiten hinsichtlich der Analyse und Verdichtung großer, heterogener Datenbestände. Für diesen Prozeß der Integration und Interpretation von Daten aus verschiedenen heterogenen Quellen sowie die darauf aufbauende Konstruktion von Modellen für einen bestimmten Problembereich mit dem Ziel der Gewinnung von Informationen einer neuen, höheren Qualität wird der Begriff *Informationsfusion* verwendet.

Aus der Zielstellung der Informationsfusion resultieren Anforderungen an Methoden und Techniken aus dem Datenbankbereich. Diese betreffen vor allem die effiziente Datenanalyse und -aufbereitung in verteilten, heterogenen Quellen und die Behandlung unzureichend strukturierter, inkonsistenter oder vager Informationen. Im vorliegenden Beitrag werden ausgehend vom aktuellen Forschungsstand wichtige Anforderungen aufgezeigt sowie potentielle Entwicklungsrichtungen und Anwendungen diskutiert.

2 Stand der Forschung

Informationsfusion ist ein interdisziplinäres Gebiet, das auf Methoden und Techniken verschiedener Bereiche, wie z.B. Datenbanken, Statistik, Maschinelles Lernen, Soft Computing oder Visualisierung zurückgreift. Nachfolgend soll der Stand der aktuellen Forschung zu dieser Thematik im wesentlichen aus Datenbanksicht skizziert werden. Die betrachteten Kerngebiete sind dabei Datenintegration und -management, Data Mining und Data Fusion.

Datenintegration und -management. In der Literatur wird Daten(bank)integration oft mit Schemaintegration gleichgesetzt. Eine ganze Reihe von Ansätzen wurden für die Schemaintegration entwickelt. Verschiedene Überblicke darüber sind z.B. in [BLN86,PBE95,Con97] zu finden. Ein zentrales Anliegen der Datenbankintegration ist die Überwindung der Heterogenität auf verschiedenen Ebenen. Im Mittelpunkt dieser Arbeiten stehen in der Regel die Heterogenitäten, die durch unterschiedliche Datenmodelle und unterschiedliche Modellierung beim Datenbankentwurf entstehen. Die

übliche Vorgehensweise bei der Integration läßt sich wie folgt (idealisiert) darstellen. Zunächst werden die zu integrierenden Datenbankschemata in ein gemeinsames Datenmodell transformiert, um die Heterogenität auf Datenmodellebene zu beseitigen. Anschließend müssen die übereinstimmenden Teile der Schemata identifiziert werden, damit dann die eigentliche Integration durchgeführt werden kann. Aufgrund der häufig anzutreffenden Heterogenität in der Modellierung desselben oder eines ähnlichen Sachverhaltes muß man sich bei diesem Schritt insbesondere mit dem Vergleich unterschiedlicher Modellierungen beschäftigen. Eine weitverbreitete Schemaarchitektur, die dieses Vorgehen unterstützt, ist die 5-Ebenen-Schema-Architektur [SL90]. Neben der Betrachtung auf Schemaebene ist für den Betrieb eines föderierten Datenbanksystems auch die Realisierung von Datenzugriffsschnittstellen zu den einzelnen Systemen wichtig. Hier gibt es bereits verschiedene Prototypsysteme, in denen Datenbankadapter entwickelt wurden (z.B. in IRO-DB [GGF⁺96]).

Für die Anfragebearbeitung in föderierten Datenbanken (siehe auch [MY95]) müssen bereits bei der Schemaintegration Abbildungsinformationen zwischen den lokalen und dem integrierten Schema festgelegt werden. Während dadurch die Anfragebearbeitung prinzipiell möglich wird, sind der Optimierung aufgrund der Heterogenität der Systeme erhebliche Grenzen gesetzt. Insbesondere müßte eine globale Instanz möglichst vollständiges Wissen über die lokalen Anfrageoptimierungsstrategien haben oder diese direkt steuern können. Aktuelle Bemühungen zu dieser Thematik sind u.a. auch auf die Optimierung von Anfragen über autonome Internet-Datenbanken (*fusion queries*) gerichtet [YPAGM98].

Für die Beantwortung von Anfragen ist die Datenqualität bei der Integration von hoher Bedeutung. Obwohl dieses Problem schon lange bekannt ist, gibt es bisher nur wenige Arbeiten, z.B. [Ger98], dazu. Dieses Problem spielt auch im Bereich des Data Warehousing [Inm96] eine erhebliche, oft aber vernachlässigte Rolle. Ein Data Warehouse kann dabei als eine materialisierte Sicht oder Integration von operativen Datenbeständen betrachtet werden, auf der z.B. entscheidungsunterstützende Auswertungen durchgeführt werden sollen, ohne das operative Geschäft zu beeinträchtigen. Ein Data Warehouse unterscheidet sich in dieser Sichtweise von einem föderierten Datenbanksystem vorrangig dadurch, daß der integrierte Datenbestand materialisiert wird, um einen effizienten Zugriff datenintensiver (OLAP-) Anwendungen zu unterstützen. Änderungen der lokalen, operativen Daten können sich so aber nicht unmittelbar auf den integrierten Datenbestand auswirken. Für Anwendungen mit der Anforderung nach möglichst aktuellen Daten oder einem erheblich anderem Zugriffsverhalten bietet sich die logische Integration in einem föderierten Datenbanksystem an [Con97].

Data Mining/KDD. Der Begriff *Knowledge Discovery in Databases (KDD)* wird im wissenschaftlichem Umfeld als „der nichttriviale Prozeß der Identifikation gültiger, neuer, potentiell nützlicher und verständlicher Muster in Datenbeständen“ [FPSS96] definiert. *Data Mining* bezeichnet in diesem Zusammenhang den Teilschritt der Suche und Bewertung von Hypothesen. Im kommerziellen Bereich wird dagegen Data Mining häufig als Synonym für KDD verwendet. KDD ist ein iterativer und interaktiver Prozeß, der die folgenden Schritte umfaßt: (1) Festlegung von Problembereich und Zielen, (2) Datensammlung und -bereinigung, (3) Auswahl und Parametrisierung der Analysefunktionen und -methoden, (4) Data Mining, (5) Bewertung und Interpretation der Ergebnisse sowie (6) Nutzung des gefundenen Wissens.

Datenbanktechnologie wird hierbei insbesondere in den Schritten 2 und 3 eingesetzt. In Abhängigkeit von der Analyseaufgabe kommen verschiedene Methoden des Data Mining zum Einsatz. Die wichtigsten Klassen dieser Verfahren sind u.a. [FPSS96,CHY96,Wro98]:

- *Erkennung von Abhängigkeiten*: Diese Verfahren ermitteln statistische Abhängigkeiten zwischen Variablen der relevanten Datensätze. Als Ergebnis werden Assoziationsregeln [AS94] oder auch Wahrscheinlichkeitsnetze geliefert.
- *Klassifikation*: Klassifikationsverfahren zielen auf die Zuordnung von Objekten zu verschiedenen vorgegebenen Klassen ab, wobei das Klassifikationsmodell anhand einer Beispielmenge (Trainingsset) der Datenbank ermittelt wird.
- *Clustering*: Beim Clustering werden ähnliche Objekte in neu gebildete Kategorien eingeordnet, so daß die Ähnlichkeiten der Objekte innerhalb einer Kategorie möglichst groß und zwischen den Kategorien gering sind [Fis95].
- *Generalisierung*: Dies beinhaltet Methoden zur Aggregation und Verallgemeinerung großer Datenmengen auf einer höheren Abstraktionsebene. Oft werden diese Verfahren bei der interaktiven Datenexploration angewendet [HCC92].
- *Sequenzanalyse*: Diese Verfahren dienen zur Suche nach häufig auftretenden Episoden oder Ereignisfolgen in Datenbeständen, denen eine (z.B. zeitliche) Ordnung der einzelnen Datensätze zugrundeliegt [MTV95].

Neben der Verarbeitung einfacher, relationaler Daten gewinnt die Analyse von Textdokumenten (*Document/Text Mining*), Bilddatenbanken (*Image Mining*), geographischen Daten (*Spatial Data Mining*) und Informationen aus dem World Wide Web (*Web Mining*) zunehmend an Bedeutung. Für eine weitergehende Diskussion konkreter KDD-Verfahren sei an dieser Stelle auf die Literatur [SHF96,HMPU97] verwiesen.

Datenfusion. Der Begriff der Datenfusion wird gegenwärtig in verschiedenen Anwendungsgebieten in teilweise unterschiedlicher Bedeutung verwendet. Grundsätzlich wird darunter jedoch die Kombination von Daten aus verschiedenen (heterogenen) Quellen verstanden. Konkrete Problemstellungen sind u.a. [LK95,PGV95,AZ98]: *Multi-Sensor Fusion*, die sich mit der Verbindung von Meßdaten verschiedener Geräte (z.B. in industriellen oder militärischen Bereichen) befaßt, *Multiple Source Interrogation* mit dem Ziel der Zusammenführung von Informationen aus Daten- oder Wissensbanken und *Image Fusion* als Kombination von Bildern einer Szene mit dem Ziel einer verbesserten Interpretation (z.B. Röntgenaufnahmen oder Satellitenbilder).

3 Anforderungen der Informationsfusion

Aufgrund der vielfältigen Problemstellungen und Anwendungsfelder der Informationsfusion und der damit verbundenen Anforderungen werden wir uns im weiteren auf die Fusionierung von Daten bzw. Informationen aus Datenbanken beschränken. Die dafür notwendigen Funktionen eines Softwaresystems lassen sich wie folgt beschreiben:

- *Datenzugriff*: Zunächst ist der transparente Zugriff auf Daten aus unterschiedlichen Quellen zu realisieren. Dies schließt die Verwendung von Datenbank-Gateways zur

Verbergung der Heterogenität ebenso ein wie die Verarbeitung von Dateien mit vorgegebener Struktur (semistrukturierte Daten), wobei der Zugriff über entsprechende Protokolle (z.B. HTTP) erfolgen kann. Weiterhin ist zu diesem Funktionsbereich die Verarbeitung und Optimierung von Anfragen zu zählen.

- *Datenintegration*: Für die Daten aus den einzelnen Quellen ist eine integrierte Sicht zu schaffen, die die Daten in einem homogenen Modell präsentiert und dabei Konflikte auf Schema- und Instanzebene behebt. Weiterhin sind quellenübergreifende Beziehungen zu repräsentieren und in geeigneter Weise zu verwalten.
- *Analyse und Verdichtung*: Durch das Extrahieren von Zusammenhängen und Abstraktionen, durch Filterung und Verdichtung der Daten sind Informationen einer neuen Qualität zu gewinnen. Die Definition der „neuen Qualität“ ist dabei abhängig von der konkreten Anwendung. Mögliche Repräsentationen für diese Informationen sind generalisierte Aggregationen und Assoziationen, Cluster und Klassen.
- *Präsentation und Weiterverarbeitung*: Die gewonnenen Informationen sind entsprechend der Problemstellung zu präsentieren bzw. zur Weiterverarbeitung bereitzustellen [KK96].
- *Repräsentation von Metainformationen*: Eine wesentliche Voraussetzung für die Fusion ist das Vorhandensein von Informationen über die Datenquellen, die zu fusionierenden Objekte und den Problembereich. Diese Metainformationen sind durch das System zu verwalten und im Verlauf des Fusionsprozesses sukzessive anzupassen bzw. zu erweitern.

Diese Funktionsbereiche sind durch eine Infrastruktur bereitzustellen, die die Basis für unterschiedliche Anwendungen der Informationsfusion bildet.

4 Anforderungen an die Datenbanktechnologie

Aus der Zielstellung der Informationsfusion und der beschriebenen Funktionalitäten leiten sich Anforderungen an Entwicklungen im Bereich Datenbanken sowie KDD ab. Für den Datenbankbereich beinhaltet dies die folgenden Aufgabenstellungen:

- *Intelligente Unterstützung des Integrationsprozesses*: Für viele Einsatzfälle ist die Integration der Schemata der einzelnen Quellen in ein globales Schema ein komplexer Prozeß, der nicht vollständig automatisierbar ist. So sind semantische und strukturelle Konflikte der Ausgangsschemata zu beseitigen und unterschiedliche Klassenhierarchien zu integrieren. Es werden daher Werkzeuge benötigt, die diese Schritte unterstützen und dabei auch die Semantik der Daten einbeziehen. Darüber hinaus sind Aspekte der Qualität der zu integrierenden Daten zu berücksichtigen, da diese die Ergebnisse der Fusion nachhaltig beeinflussen.
- *Realisierung eines effizienten Datenzugriffs*: Zur Analyse großer Datenbestände über verschiedene Quellen hinweg sind effiziente Zugriffsmechanismen notwendig. Speziell für verteilte, heterogene Quellen sind geeignete Indexstrukturen aufzubauen und spezielle Caching- oder Replikationsstrategien zu verfolgen. Weiterhin sind die spezifischen Anforderungen der Analysemethoden hinsichtlich der Zugriffsschnittstellen (z.B. satzorientierte oder navigierende Zugriffe) zu berücksichtigen.

- *Integration semistrukturierter Daten:* Nicht zuletzt durch die Verbreitung des World Wide Web liegen viele Informationen in nur unzureichend strukturierter Form, wie z.B. in HTML-Dateien, vor. Die Integration dieser semistrukturierten Daten, die effiziente Aufbereitung und Repräsentation sowie die damit verbundene Möglichkeit der Anfragebearbeitung stellt eine weitere wichtige Aufgabe dar.
- *Gewinnung von Metainformationen:* Informationen, die die Semantik und Qualität der Daten beschreiben, bilden eine wichtige Basis für die Fusion. Sofern diese Metainformationen nicht vorliegen, müssen sie aus den Daten extrahiert oder mit Hilfe des Nutzers erfaßt werden.

Ein Teil dieser Themen sind aktueller Forschungsgegenstand z.B. im Umfeld föderierter Datenbanken. Aus den aufgeführten Aufgabenstellungen lassen sich konkrete Anforderungen an Datenbankmanagementsysteme (DBMS) ableiten:

- Das eingesetzte DBMS muß einen offenen Optimierer haben, um Fusionsmethoden zusammen mit DB-Operationen optimieren zu können. So werden in Fusionsprozessen in der Regel (statistische) Aufbereitungsschritte im Wechsel mit Filterungsschritten zur Selektion eingesetzt, die durch DBMS-Anfragen effizient unterstützt werden können. Hier muß eine übergreifende Optimierung erfolgen.
- Das DBMS muß Funktionalität zur Unterstützung der Integration externer Daten anbieten, die über einfache Import/Export-Routinen hinausgeht.
- Das DBMS sollte Funktionen zur Unterstützung des Rankings bzw. der Qualitätsbewertung von Anfrageergebnissen besitzen (also müssen in gewissem Umfang Techniken des Information Retrieval in das DBMS verlagert werden).
- Das DBMS muß eine offene Softwarearchitektur mit einem zugänglichen Repository für Metainformationen besitzen, um Fusionsmethoden einbetten zu können.
- Für die (physische) Anfragebearbeitung sollten neue Techniken, wie etwa die Indexgenerierung *on the fly*, aber auch die Nutzung und Integration spezieller vorhandener Indexstrukturen in den verschiedenen Datenquellen verfügbar sein.
- Viele statistische Methoden benutzen eine zufällige Auswahl von Datensätzen als ersten Schritt zur Initialisierung, bevor der gesamte Datenbestand analysiert wird. Ein derartiges *Sampling* zur Generierung einer zufälligen Stichprobe wird von kommerziellen DBMS in der Regel nicht unterstützt.

5 Anforderungen an den KDD-Prozeß

Im Bereich KDD betreffen die Anforderungen die folgenden Schwerpunkte:

- *Behandlung verschiedener Datentypen:* Neben einfachen relationalen Datentypen sind auch komplexe Typen im Analyseprozeß zu berücksichtigen. Hierzu gehört die Verarbeitung strukturierter Daten und Objekte, von Hypertext- und Multimediatdaten oder auch von Daten mit einem zeitlichen oder räumlichen Bezug.
- *Behandlung unsicherer und vager Daten:* Bei der Fusion sind oft auch unvollständige, ungenaue, vage oder auch Erfahrungsdaten zu verarbeiten. Diese Unsicherheit/Vagheit ist zunächst geeignet zu modellieren (z.B. mit Wahrscheinlichkeitsverteilungen, mehrwertigen Logiken, Fuzzy-Mengen u.ä.) und in den Fusionsprozeß einzubeziehen.

- *Effizienz und Skalierbarkeit der Verfahren:* Für die Analyse sehr großer Datenbestände sind Effizienz und Skalierbarkeit von großer Bedeutung. Hauptspeicherbasierte Verfahren, die bei Datenmengen im Megabyte-Bereich akzeptabel arbeiten, können im Gigabyte-Bereich aufgrund von Speichermangel versagen. Außerdem sollte die Laufzeit der Verfahren in Abhängigkeit von der Datenmenge vorhersehbar bzw. abschätzbar sein.
- *Verbesserung der Aussagefähigkeit und Verständlichkeit der Ergebnisse:* Der Nutzen der Fusion für ein gegebenes Problem hängt entscheidend von der Aussagekraft und Überschaubarkeit der Ergebnisse ab. Hierzu sind zunächst Rauschen, Ausnahmen oder irrelevante Zusammenhänge zu unterdrücken. Weiterhin sind die für das Analyseproblem geeigneten Ausdrucksformen zu wählen (z.B. Regeln, Entscheidungsbäume, graphische Darstellungen etc.) und die Qualität der Ergebnisse auch unter Berücksichtigung der Qualität der Eingangsdaten anzugeben.

6 Anforderungen an Systemarchitektur

Aus der Gesamtsicht eines Systems zur Informationsfusion stellen sich darüber hinaus noch folgende Forderungen:

- *Unterstützung einer interaktiven und iterativen Arbeitsweise:* Aufgrund der Komplexität des Fusionsprozesses und des Volumens der zu verarbeitenden Daten ist eine schrittweise, interaktive Durchführung der Fusion notwendig. So kann zunächst mit wenigen, ausgewählten Daten eines relevanten Teilschemas die Anwendbarkeit bestimmter Methoden und die zu erwartende Qualität der Ergebnisse mit geringem Aufwand abgeschätzt werden, bevor die Fusion über den gesamten integrierten Datenbestand erfolgt. Weiterhin ist eine schrittweise Fokussierung auf bestimmte Datenbereiche wünschenswert. Aus der Beobachtung der Kosten für Datenzugriff und -analyse ergibt sich außerdem die Möglichkeit einer Optimierung des Laufzeitverhaltens durch geeignete Caching-Strategien, die Verwendung speziell angepaßter Indexstrukturen, die Parallelisierung der Verarbeitung oder auch durch die Umordnung von Analyseoperationen.
- *Anpassungsfähigkeit und Erweiterbarkeit des Systems:* Auch wenn ein Fusionssystem für einen konkreten Anwendungsfall entwickelt wird, können nicht immer die zu unterstützenden Datenquellen und die benötigten Methoden vorausbestimmt werden. Die Erweiterbarkeit des Fusionssystems um neue Methoden (z.B. durch Plugins [WWSE96]) und Datenquellen ist daher ein weiteres wichtiges Kriterium.
- *Intelligente Nutzerunterstützung bei Auswahl und Anwendung der Fusionsmethoden:* Der Einsatz der verschiedenen Integrations-, Fusions- und Analysemethoden setzt tiefe Kenntnisse über den Problembereich, die Struktur und Semantik der Daten sowie die Methoden selbst voraus. Die Auswahl, Kombination und gegebenenfalls Parametrisierung der Methoden ist durch geeignete Techniken (z.B. Nutzung von Metainformationen und Wissen zum Problembereich, Voranalyse der Daten, Einsatzkriterien) zu unterstützen.

Mit den heute verfügbaren Techniken sind einige dieser Anforderungen bereits erfüllbar. Forschungsbedarf besteht aber noch bei der Verbindung dieser Techniken, wie z.B.

der Fusion von Daten unterschiedlicher Qualität und Repräsentation, der effizienten Analyse verteilter, heterogener Quellen und der kombinierten Anwendung verschiedener Analyse- und Fusionsmethoden.

7 Anwendungsfelder

Potentielle Anwendungen der Informationsfusion existieren überall dort, wo Daten aus verschiedenen Quellen zu kombinieren sind, um daraus neue Informationen abzuleiten und Entscheidungsprozesse zu unterstützen. Es bestehen so durchaus Anknüpfungspunkte zum Data Warehousing, wobei jedoch Informationsfusion auf die intelligente (semi-)automatische Transformation großer heterogener Datenmengen ausgerichtet ist, während beim Data Warehousing die interaktive Exploration von im Vorfeld integrieren und redundant gehaltenen Beständen im Vordergrund steht. Zwei weitere Anwendungsszenarien der Informationsfusion sollen im folgenden kurz skizziert werden. Die beiden Bereiche werden in der Magdeburger Datenbankgruppe gegenwärtig in Kooperationsprojekten mit Fachexperten bearbeitet.

Bioinformatik. Im Rahmen von Forschungsbemühungen auf dem Gebiet der Biotechnologien werden weltweit Daten mit molekulargenetischen Zusammenhängen gesammelt und zum Teil über das Internet verfügbar gemacht. Diese Datenbestände, die meist in Form einfacher Dateien mit proprietären Formaten vorliegen, repräsentieren aber oft nur ausgewählte Aspekte einer konkreten Anwendung (wie Gene, Informationen zu Krankheiten oder Stoffwechselkreisläufe). Durch eine Integration dieser Datenbestände [HHSS98] und die darauf aufbauende Analyse mit KDD-Techniken können neue Zusammenhänge aufgedeckt werden. Erschwert wird dies jedoch durch die Struktur, Inkonsistenz und Fehlerbehaftung der Daten sowie die ständige Aktualisierung und Erweiterung, so daß hier ein wichtiges Anwendungsfeld für Fusionstechniken besteht.

Telekommunikation. Ein weiteres Szenario aus dem Bereich der Telekommunikation wird in [Mat97] als Data Warehouse-Anwendung vorgestellt. Die Betreiber großer Telekommunikationsanlagen verwalten in ihren einzelnen Geschäftsbereichen eine Vielzahl Daten zu Produkten, Kunden und ihren Verbindungen sowie zum Telekommunikationsnetz selbst. Die Integration und Verdichtung dieser Datenbestände ist nicht nur für die dispositiven Bereiche oder das Marketing von großer Bedeutung, sondern kann auch das Netzwerk- und Systemmanagement wirkungsvoll unterstützen. So ist eine Teilaufgabe des Fehlermanagements die Alarmkorrelation: die Zuordnung von Fehlerursache und -behandlung zu Störungen des Netzwerkes, die durch Alarme repräsentiert werden [JW93]. Die Definition des notwendigen Korrelationsmodells ist insbesondere für große Netze sehr komplex. Aus der Analyse von historischen Alarmsequenzen lassen sich jedoch Regeln zur Generalisierung und Komprimierung von Alarmepisoden ableiten [HKM⁺96], die in Verbindung mit Informationen aus der Konfigurationsdatenbank die Basis des Korrelationsmodells bilden können.

8 Ausblick

Die Kombination von Daten aus unterschiedlichen Quellen sowie die darauf aufbauende Gewinnung von neuen Informationen durch Filterung, Verdichtung und Extraktion von

Zusammenhängen ist eine Aufgabe, die in vielen Anwendungsgebieten besteht. Insbesondere auch zur Beherrschung des stetig zunehmenden Informationsaufkommens aufgrund des einfachen Zugriffs auf weltweit verfügbare Quellen wächst der Bedarf nach einer intelligenten Informationsfusion. Einen wichtigen Beitrag hierzu müssen Methoden und Techniken aus dem Datenbank- und Data Mining-Bereich leisten.

Im Rahmen eines geplanten Forschungsvorhabens, an dem mehrere Forschungsgruppen der Universität Magdeburg beteiligt sind, sollen wichtige Aspekte der beschriebenen Anforderungen hinsichtlich der Verbindung von Integrations-, Fusions- und Analysemethoden für heterogene, verteilte Datenbestände untersucht werden. Hierbei wird insbesondere angestrebt, das Fachwissen aus verschiedenen Bereichen gezielt zusammenzubringen, um so einerseits anwendungsspezifisches Wissen berücksichtigen zu können und andererseits die Anforderungen aus den verschiedenen beteiligten Bereichen effektiv und möglichst effizient erfüllen zu können.

Danksagung

Die Autoren danken R. Kruse und den anderen Magdeburger Kollegen für die ausführliche und fruchtbare Diskussion über Informationsfusion.

Literatur

- [AS94] R. Agrawal und R. Srikant. Fast Algorithms for Mining Association Rules in Large Databases. In *Proc. of the 20th Int. Conf. on Very Large Data Bases (VLDB)*, S. 478–499, Santiago, Chile, September 1994.
- [AZ98] H. Arabnia und D. Zhu, Herausgeber. *Proc. of the Int. Conf. on Multisource-Multisensor Information Fusion - FUSION '98*, Las Vegas, NV, 1998. CSREA Press.
- [BLN86] C. Batini, M. Lenzerini und S. Navathe. A Comparative Analysis of Methodologies for Database Schema Integration. *ACM Computing Surveys*, 18(4):323–364, 1986.
- [CHY96] M. Chen, J. Han und P. Yu. Data Mining: An Overview from a Database Perspective. *IEEE Transactions on Knowledge and Data Engineering*, 8(6):866–883, 1996.
- [Con97] S. Conrad. *Föderierte Datenbanksysteme: Konzepte der Datenintegration*. Springer-Verlag, Berlin/Heidelberg, 1997.
- [Fis95] D. Fisher. Optimization and simplification of hierarchical clustering. In *Proc. of 1st Int. Conf. on Knowledge Discovery and Data Mining (KDD-95)*, S. 118–123, Montreal, Canada, August 1995.
- [FPSS96] U. Fayyad, G. Piatetsky-Shapiro und P. Smyth. From Data Mining to Knowledge Discovery: An Overview. In U. Fayyad, G. Piatetsky-Shapiro, P. Smyth und R. Uthurusamy, Herausgeber, *Advances in Knowledge Discovery and Data Mining*, Kapitel 1, S. 1–34. AAAI/MIT Press, Cambridge, MA, 1996.
- [Ger98] M. Gertz. Managing Data Quality and Integrity in Federated Databases. In *2nd Annual IFIP TC-11 WG 11.5 Working Conf. on Integrity and Internal Control in Information Systems*, Warrenton, Virginia, November 1998. To appear.
- [GGF⁺96] G. Gardarin, S. Gannouni, B. Finance, P. Fankhauser, W. Klas, D. Pastre, R. Legoff und A. Ramfos. IRO-DB — A Distributed System Federating Object and Relational Databases. In *Object-Oriented Multidatabase Systems — A Solution for Advanced Applications*, Kapitel 20, S. 684–712. Prentice Hall, Eaglewoods Cliffs, NJ, 1996.

- [HCC92] J. Han, Y. Cai und N. Cercone. Knowledge Discovery in Databases: An Attribute-Oriented Approach. In *Proc. of 1992 Int. Conf. on Very Large Data Bases (VLDB'92)*, S. 547–559, Vancouver, Canada, August 1992.
- [HHSS98] M. Höding, R. Hofestädt, G. Saake und U. Scholz. Schema Derivation for WWW Information Sources and their Integration with Databases in Bioinformatics. In *Advances in Databases and Information Systems – ADBIS'98, Poznań, Poland, September 1998*, LNCS 1475, S. 296–304, Berlin, 1998. Springer-Verlag.
- [HKM⁺96] K. Hätönen, M. Klemettinen, H. Mannila, P. Ronkainen und H. Toivonen. Knowledge Discovery from Telecommunication Network Alarm Databases. In *Proc. of 12th Int. Conf. on Data Engineering (ICDE'96)*, S. 115–122, New Orleans, 1996.
- [HMPU97] D. Heckerman, H. Mannila, D. Pregibon und R. Uthurusamy, Herausgeber. *KDD-97 – Proc. of the 3rd Int. Conf. on Knowledge Discovery and Data Mining*, Menlo Park, CA, 1997. AAAI Press.
- [Inm96] W. H. Inmon. *Building the Data Warehouse*. Wiley & Sons, 2 Auflage, 1996.
- [JW93] G. Jakobson und M.D. Weissman. Alarm Correlation. *IEEE Network*, 7(6):52–59, November 1993.
- [KK96] D. Keim und H.-P. Kriegel. Visualization Techniques for Mining Large Databases: A Comparison. *IEEE Transactions on Knowledge and Data Engineering*, 8(6):923–938, December 1996.
- [LK95] R.C. Luo und M.G. Kay, Herausgeber. *Multisensor Integration and Fusion for Intelligent Machines and Systems*. Ablex Publishing Corporation, Norwood, NJ, 1995.
- [Mat97] R. Mattison. *Data Warehousing and Data Mining for Telecommunications*. Artech House, Norwood, MA, 1997.
- [MTV95] H. Mannila, H. Toivonen und A.I. Verkano. Discovering frequent episodes in sequences. In *Proc. of 1st Int. Conf. on Knowledge Discovery and Data Mining (KDD-95)*, S. 210–215, Montreal, Canada, August 1995.
- [MY95] W. Meng und C. Yu. Query Processing in Multidatabase Systems. In W. Kim, Herausgeber, *Modern Database Systems*, S. 551–572, New York, NJ, 1995. ACM Press.
- [PBE95] E. Pitoura, O. Bukhres und A. K. Elmagarmid. Object Orientation in Multidatabase Systems. *ACM Computing Surveys*, 27(2):141–195, 1995.
- [PGV95] S. Pfeleger, J. Goncalves und D. Vernon, Herausgeber. *Data Fusion Applications*. Springer-Verlag, Berlin, 1995. Research Reports ESPRIT.
- [SHF96] E. Simoudis, J. Han und U. Fayyad, Herausgeber. *KDD-96 – Proc. of the 2nd Int. Conference on Knowledge Discovery and Data Mining*, Menlo Park, CA, 1996. AAAI Press.
- [SL90] A. P. Sheth und J. A. Larson. Federated Database Systems for Managing Distributed, Heterogeneous, and Autonomous Databases. *ACM Computing Surveys*, 22(3):183–236, 1990.
- [Wro98] S. Wrobel. Data Mining und Wissensentdeckung in Datenbanken. *Künstliche Intelligenz – Organ des FB I der Gesellschaft für Informatik (GI)*, (1), 1998.
- [WWSE96] S. Wrobel, D. Wettschereck, E. Sommer und W. Emde. Extensibility in data mining systems. In Simoudis et al. [SHF96].
- [YPAGM98] R. Yerneni, Y. Papakonstantinou, S. Abiteboul und H. Garcia-Molina. Fusion queries over internet databases. In *Advances in Database Technology - EDBT'98*, LNCS 1377, S. 57–71. Springer-Verlag, 1998.