



Informationsintegration

Anwendungsszenarien

Ulf Leser

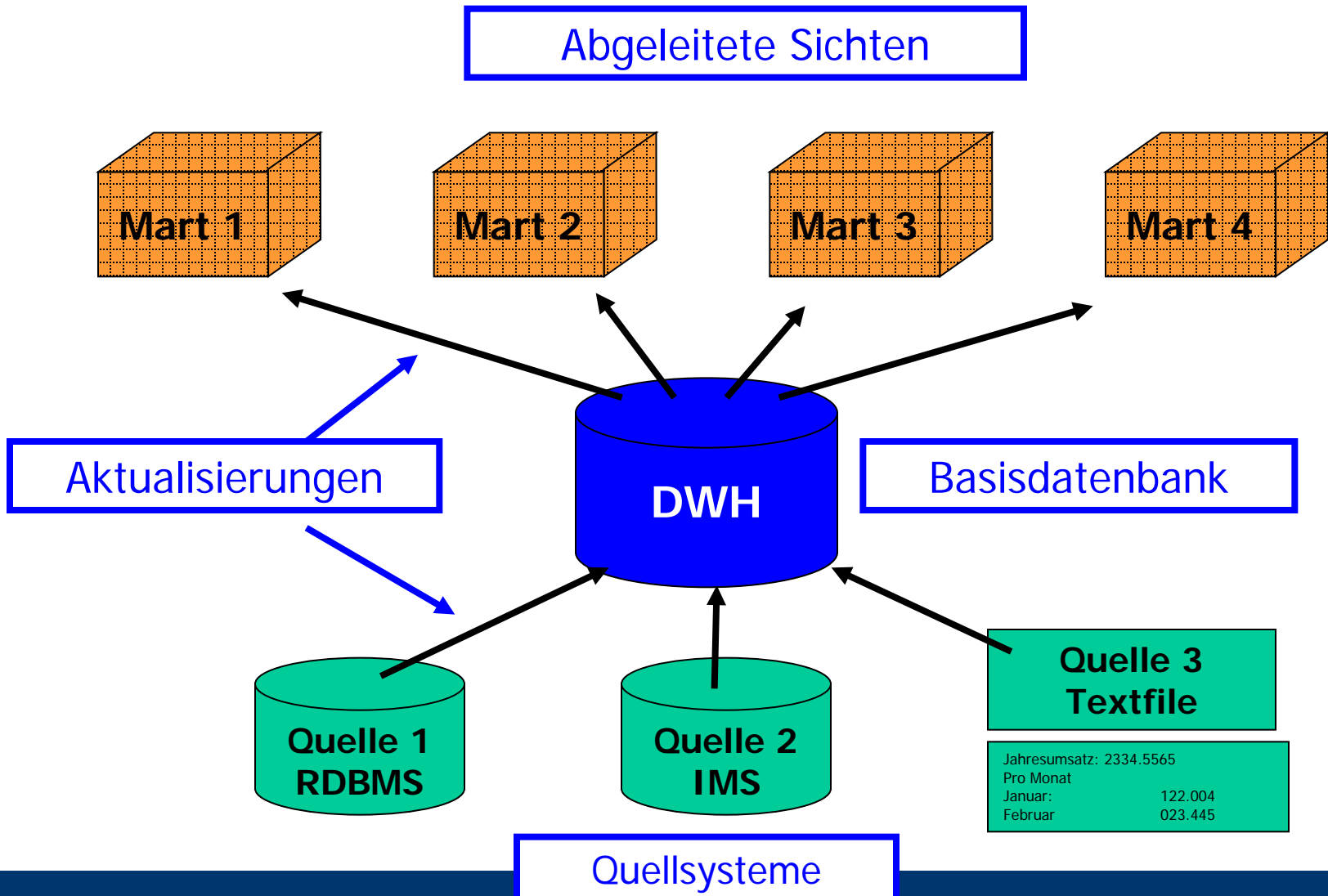
Inhalt dieser Vorlesung

- Zwei Anwendungsszenarien
 - Data Warehouse
 - Föderierte Datenbanken
- Abgrenzung
- Daten versus Schema

Data Warehouse (siehe eigene Vorlesung)

- Integrierte Datenbank innerhalb eines Unternehmens
- Ziel: Zugriff auf alle Informationen, um **strategische Entscheidungen** zu unterstützen
- Charakteristisch: Aufbau einer **integrierten Datenbasis**
 - Daten werden aus Quellen kopiert
 - Zur Anfragezeit wird nur das Data Warehouse benutzt
 - Problem: Aktualisierung
- Homogenisierung beim Import
 - Einheiten, Attributnamen, Formate, ...
- Integrationsprozess im DWH heißt **ETL**
 - Extraction, Transformation, Load

Übersicht



Frage eines Biologen

Finde alle menschlichen Sequenzen, die zu mindestens 60% identisch sind mit „channel“ Proteinen, die im Gewebe des zentralen Nervensystems der Maus exprimiert werden



Quelle: *A Practitioner's Guide to Data Management and Data Integration in Bioinformatics*, Barbara A. Eckman in Bioinformatics by Zoe Lacroix and Terence Critchlow, 2003, Morgan Kaufmann.

Verschiedene Informationsquellen

- Beteiligte Informationsquellen
 - Mouse Genome Database (MGD) @ Jackson Labs
 - SwissProt @ EBI
 - BLAST tool @ NCBI
 - GenBank nucleotide sequence database @ NCBI



- Alle Quellen sind frei verfügbar

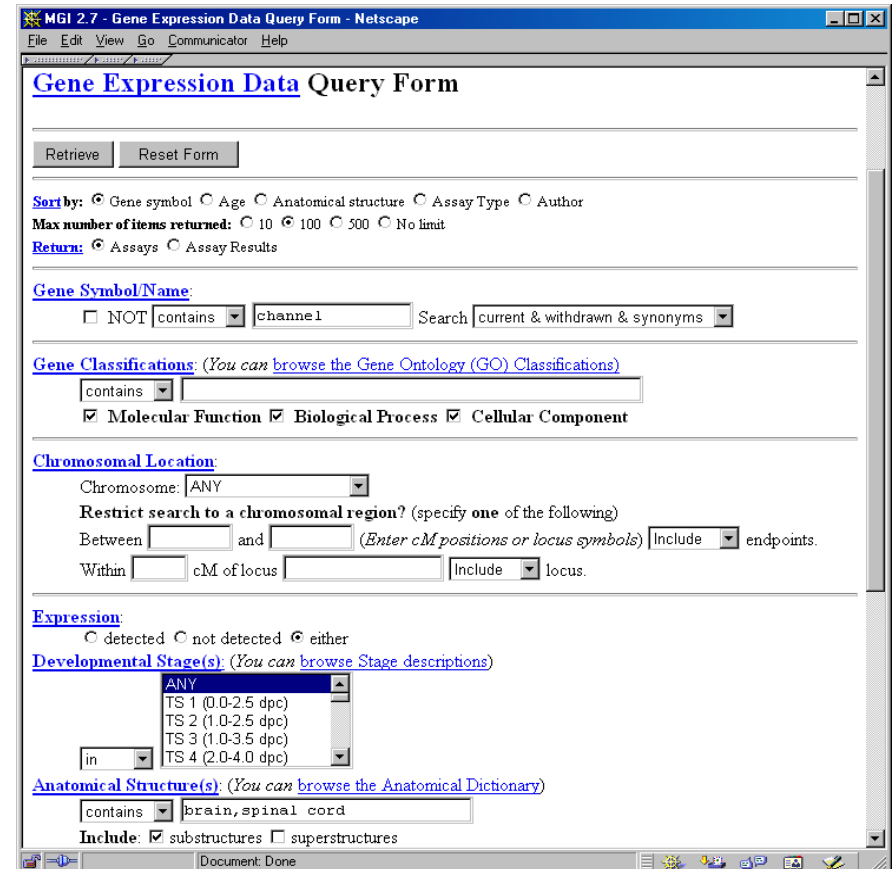
Zusammenhang

Finde alle menschlichen Sequenzen, die zu mindestens 60% identisch sind mit „channel“ Proteinen, die im Gewebe des zentralen Nervensystems der Maus exprimiert werden

- MGD speichert Aktivität und Sequenz von Mausgenen
 - Was sind „channel“ Proteine? Wo werden die exprimiert?
- Swiss-Prot speichert die Proteinsequenz eines Gens
 - Proteinsequenz zu exprimierten Genen
- BLAST findet ähnliche Proteinsequenzen
 - Suche nach ähnlichen Sequenzen
- Genbank enthält DNA-Sequenz von Proteinen
 - Spezies zu einem Protein? DNA-Sequenz zum Protein
- **Keine Datenbank enthält alles**

Herkömmlicher Ansatz: Browsing

1. Suche im zentralen Nervensystem aktive „channel“ Gene und deren Sequenzen in der MGD mittels HTML Formular



The screenshot shows the MGI 2.7 Gene Expression Data Query Form in a Netscape browser window. The form is titled "Gene Expression Data Query Form" and includes several sections for filtering search results:

- Retrieve/Reset Form:** Buttons for "Retrieve" and "Reset Form".
- Sort by:** Radio buttons for "Gene symbol", "Age", "Anatomical structure", "Assay Type", and "Author".
- Max number of items returned:** Radio buttons for "10", "100", "500", and "No limit".
- Return:** Radio buttons for "Assays" and "Assay Results".
- Gene Symbol/Name:** A search box containing "channel" and a dropdown menu set to "current & withdrawn & synonyms".
- Gene Classifications:** A dropdown menu set to "contains" and three checked checkboxes: "Molecular Function", "Biological Process", and "Cellular Component".
- Chromosomal Location:** A dropdown menu set to "ANY", a section for "Restrict search to a chromosomal region?" with "Between" and "Within" options, and "Include" dropdown menus for "endpoints" and "locus".
- Expression:** Radio buttons for "detected", "not detected", and "either".
- Developmental Stage(s):** A dropdown menu set to "ANY" with a list of stages: "TS 1 (0.0-2.5 dpc)", "TS 2 (1.0-2.5 dpc)", "TS 3 (1.0-3.5 dpc)", and "TS 4 (2.0-4.0 dpc)".
- Anatomical Structure(s):** A dropdown menu set to "contains" and a text box containing "brain,spinal cord".
- Include:** Checked checkboxes for "substructures" and "superstructures".

Herkömmlicher Ansatz: Browsing

- MGD Resultat
 - 14 Gene aus 17 Experimenten

MGI 2.7 - Gene Expression Data Query Results (Summary) - Netscape

File Edit View Go Communicator Help

[Gene Expression Data](#)

Query Results -- Summary

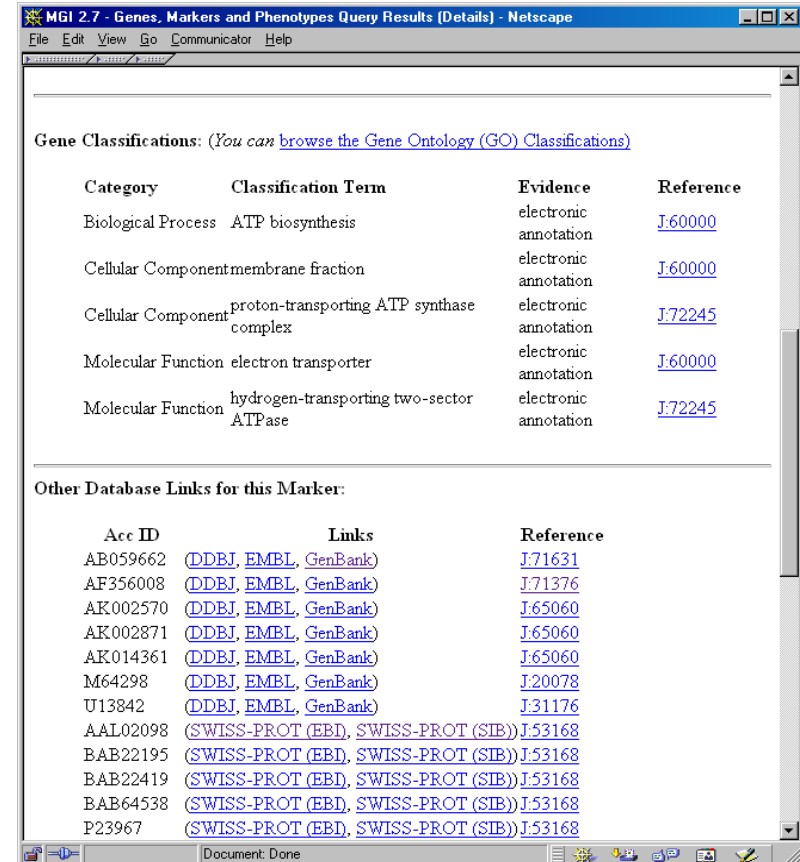
17 matching assays displayed

Gene	Assay Type	Assay	RefID	Reference
Atp6l	Northern blot	MGI2150866	J71376	Nishi T, J Biol Chem 2001 Sep 7;276(36):34122-30
Cacnb3	RT-PCR	MGI1205020	J46439	Freeman TC, MGI Direct Data Submission 1998,0
Gja1	Immunohistochemistry	MGI1338492	J31725	Yancey SB, Development 1992 Jan;114(1):203-12
Gja1	Immunohistochemistry	MGI1338557	J31725	Yancey SB, Development 1992 Jan;114(1):203-12
Kcna4	Immunohistochemistry	MGI1335744	J41027	Zhong W, Development 1997 May;124(10):1887-97
Kcnab2	RT-PCR	MGI1204928	J46439	Freeman TC, MGI Direct Data Submission 1998,0
Kcnh1	RT-PCR	MGI1205795	J46439	Freeman TC, MGI Direct Data Submission 1998,0
Kcnj12	RT-PCR	MGI1204727	J46439	Freeman TC, MGI Direct Data Submission 1998,0
Kcnj2	RT-PCR	MGI1205781	J46439	Freeman TC, MGI Direct Data Submission 1998,0
Kcnj3	RT-PCR	MGI1205497	J46439	Freeman TC, MGI Direct Data Submission 1998,0
Kcnj4	RT-PCR	MGI1204196	J46439	Freeman TC, MGI Direct Data Submission 1998,0
Kcnj4	RT-PCR	MGI1204198	J46439	Freeman TC, MGI Direct Data Submission 1998,0
Kcnj5	RT-PCR	MGI1205098	J46439	Freeman TC, MGI Direct Data Submission 1998,0
Kcnj6	RT-PCR	MGI1204201	J46439	Freeman TC, MGI Direct Data Submission 1998,0
Kcnj9	RT-PCR	MGI1204204	J46439	Freeman TC, MGI Direct Data Submission 1998,0
Kcnma1	RT-PCR	MGI1205940	J46439	Freeman TC, MGI Direct Data Submission 1998,0
Kcnma1	RT-PCR	MGI1205942	J46439	Freeman TC, MGI Direct Data Submission 1998,0

Document: Done

Herkömmlicher Ansatz: Browsing

- In MGD Details zu jedem der 14 Gene ansehen
- Durchschnittlich fünf SwissProt Links pro Gen



Gene Classifications: (You can [browse the Gene Ontology \(GO\) Classifications](#))

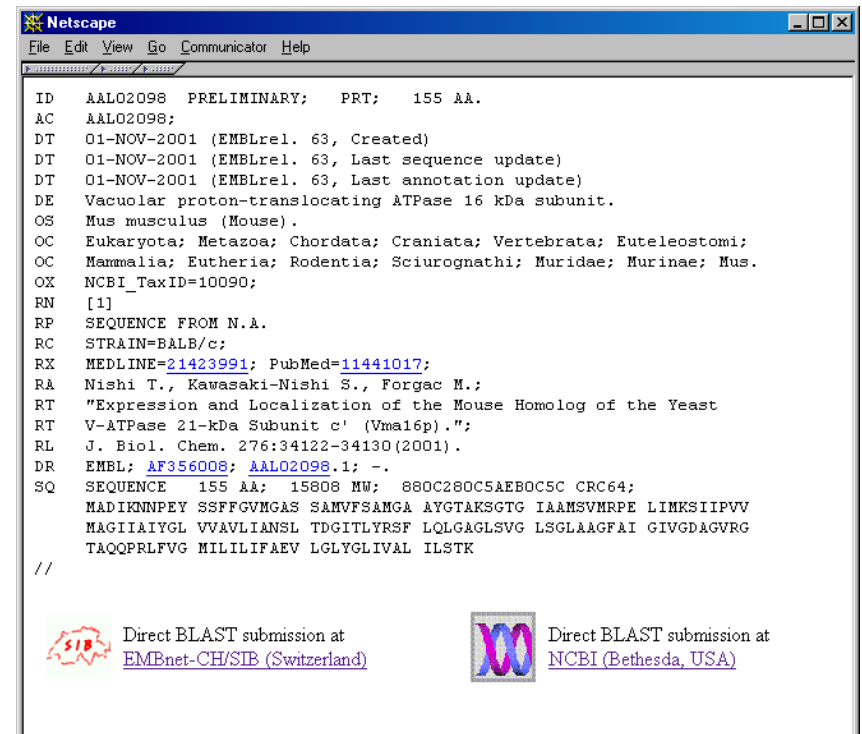
Category	Classification Term	Evidence	Reference
Biological Process	ATP biosynthesis	electronic annotation	J:60000
Cellular Component	membrane fraction	electronic annotation	J:60000
Cellular Component	proton-transporting ATP synthase complex	electronic annotation	J:72245
Molecular Function	electron transporter	electronic annotation	J:60000
Molecular Function	hydrogen-transporting two-sector ATPase	electronic annotation	J:72245

Other Database Links for this Marker:


Acc ID	Links	Reference
AB059662	(DDBJ, EMBL, GenBank)	J:71631
AF356008	(DDBJ, EMBL, GenBank)	J:71376
AK002570	(DDBJ, EMBL, GenBank)	J:65060
AK002871	(DDBJ, EMBL, GenBank)	J:65060
AK014361	(DDBJ, EMBL, GenBank)	J:65060
M64298	(DDBJ, EMBL, GenBank)	J:20078
U13842	(DDBJ, EMBL, GenBank)	J:31176
AAL02098	(SWISS-PROT (EBI), SWISS-PROT (SIB))	J:53168
BAB22195	(SWISS-PROT (EBI), SWISS-PROT (SIB))	J:53168
BAB22419	(SWISS-PROT (EBI), SWISS-PROT (SIB))	J:53168
BAB64538	(SWISS-PROT (EBI), SWISS-PROT (SIB))	J:53168
P23967	(SWISS-PROT (EBI), SWISS-PROT (SIB))	J:53168


Herkömmlicher Ansatz: Browsing

- Betrachtung jedes SwissProt Eintrags
- Durch Klick BLAST anwerfen



```
ID AAL02098 PRELIMINARY; PRT; 155 AA.
AC AAL02098;
DT 01-NOV-2001 (EMBLrel. 63, Created)
DT 01-NOV-2001 (EMBLrel. 63, Last sequence update)
DT 01-NOV-2001 (EMBLrel. 63, Last annotation update)
DE Vacuolar proton-translocating ATPase 16 kDa subunit.
OS Mus musculus (Mouse).
OC Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
OC Mammalia; Eutheria; Rodentia; Sciurognathi; Muridae; Murinae; Mus.
OX NCBI_TaxID=10090;
RN [1]
RP SEQUENCE FROM N.A.
RC STRAIN=BALB/c;
RX MEDLINE=21423991; PubMed=1441017;
RA Nishi T., Kawasaki-Nishi S., Forgac M.;
RT "Expression and Localization of the Mouse Homolog of the Yeast
RT V-ATPase 21-kDa Subunit c' (Vma16p).";
RL J. Biol. Chem. 276:34122-34130(2001).
DR EMBL; AF356008; AAL02098.1; -.
SQ SEQUENCE 155 AA; 15808 MW; 880C280C5AEB0C5C CRC64;
MADIKNNPEY SSFFGVMGAS SAMVFSAMGA AYGTAKSGTG IAAMSVMRPE LINKSIIPVV
MAGIIATYGL VVAVLIANSL TDGITLYRSF LQLGAGLSVG LSGLAAGFAT GIVGDAGVRG
TAQQPRLFVG MILILIFAEV LGLYGLIVAL ILSTK
//
```

 Direct BLAST submission at
[EMBNet-CH/SIB \(Switzerland\)](#)

 Direct BLAST submission at
[NCBI \(Bethesda, USA\)](#)

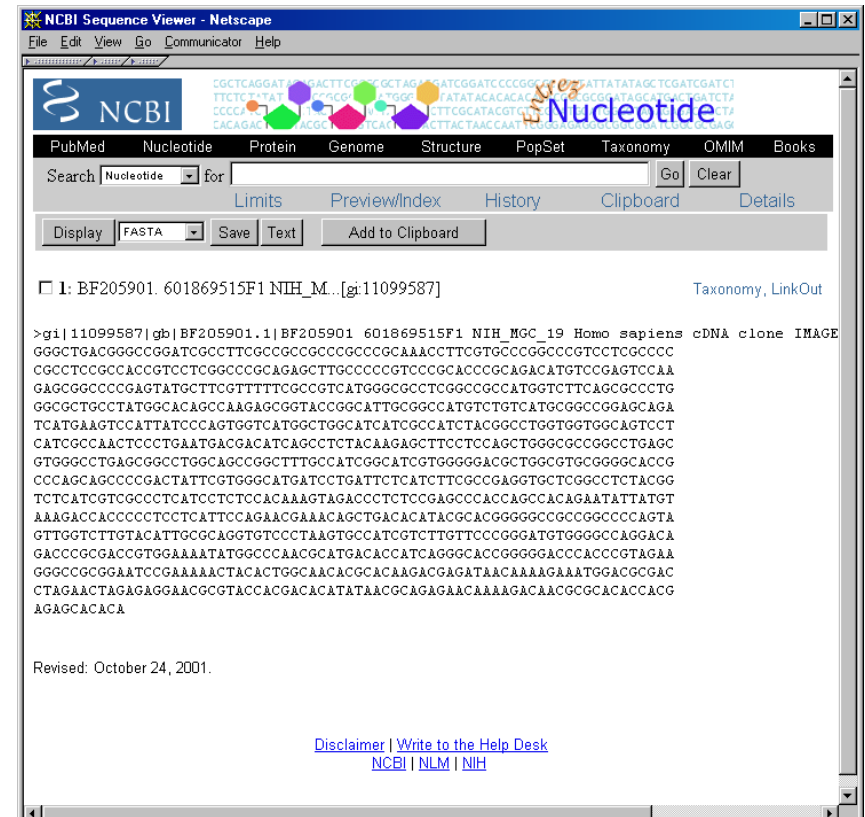
Herkömmlicher Ansatz: Browsing

- Betrachtung jedes BLAST Resultats
 - nicht-menschliche Treffer eliminieren,
 - 60% Identität sicherstellen

Sequences producing significant alignments:	Score (bits)	E Value
gb BF383501.1 BF383501 602045186F1 NCI_CGAP_Li9 Mus musculus cDN...	251	8e-67
gb BI555191.1 BI555191 603236136F1 NIH_CGAP_Man3 Mus musculus cD...	251	8e-67
gb BE285425.1 BE285425 601096726F1 NCI_CGAP_Man5 Mus musculus cD...	251	8e-67
gb BI651120.1 BI651120 603297590F1 NIH_CGAP_Man3 Mus musculus cD...	251	8e-67
...		
gb AW921969.1 AW921969 EST353273 Rat gene index, normalized rat,...	251	8e-67
gb BI646796.1 BI646796 603276734F1 NIH_CGAP_Man3 Mus musculus cD...	251	8e-67
gb BF123349.1 BF123349 601759145F1 NCI_CGAP_Man5 Mus musculus cD...	251	8e-67
gb BI693533.1 BI693533 603341913F1 NCI_CGAP_Man2 Mus musculus cD...	251	8e-67
gb BI666010.1 BI666010 603287067F1 NCI_CGAP_Man6 Mus musculus cD...	251	8e-67
...		
emb AL633622.1 AL633622 AL633622 XGC-gastrula Silurana tropicali...	228	7e-60
emb AL639595.1 AL639595 AL639595 XGC-neurola Silurana tropicalis...	228	7e-60
emb AL594253.1 AL594253 AL594253 XGC-gastrula Silurana tropicali...	228	7e-60
gb AL557998.1 AL557998 AL557998 LTI_NFL008_TC2 Homo sapiens cDN...	227	2e-59
gb BE397494.1 BE397494 601288884F1 NIH_MGC_8 Homo sapiens cDNA c...	227	2e-59
gb BF205901.1 BF205901 601869515F1 NIH_MGC_19 Homo sapiens cDNA ...	227	2e-59
gb BE729329.1 BE729329 601561519F1 NIH_MGC_20 Homo sapiens cDNA ...	227	2e-59
gb BG697408.1 BG697408 602661186F1 NCI_CGAP_Skn3 Homo sapiens cD...	227	2e-59
gb BI765781.1 BI765781 603046569F1 NIH_MGC_116 Homo sapiens cDNA...	227	2e-59
gb BE797916.1 BE797916 601586263F1 NIH_MGC_7 Homo sapiens cDNA c...	227	2e-59
gb BG490168.1 BG490168 602519116F1 NIH_MGC_18 Homo sapiens cDNA ...	227	2e-59
gb BG741416.1 BG741416 602631991F1 NCI_CGAP_Skn3 Homo sapiens cD...	227	2e-59
gb AI114460.1 AI114460 HA1042 Human fetal liver cDNA library Hom...	227	2e-59
gb AW249148.1 AW249148 2820881.Sprime NIH_MGC_7 Homo sapiens cDN...	227	2e-59
gb BF727320.1 BF727320 by19h07.y1 Human Lens cDNA (Un-normalized...	227	2e-59
gb BI328911.1 BI328911 602980634F1 NCI_CGAP_Li9 Mus musculus cDN...	226	3e-59
gb BF101272.1 BF101272 601754562F1 NCI_CGAP_Mam1 Mus musculus cD...	192	4e-59
emb AL627969.1 AL627969 AL627969 XGC-gastrula Silurana tropicali...	225	6e-59
emb AL643955.1 AL643955 AL643955 XGC-neurola Silurana tropicalis...	225	6e-59
gb BI706803.1 BI706803 fq10e10.y1 Zebrafish adult retina cDNA Pa...	225	8e-59
gb BE789647.1 BE789647 601481404F1 NIH_MGC_68 Homo sapiens cDNA ...	225	8e-59
gb BI447381.1 BI447381 dah87e12.y1 NICHD XGC Emb2 Xenopus laevis...	225	8e-59
gb BI475274.1 BI475274 fq30d06.y3 zebrafish adult brain Danio re...	225	8e-59
gb AW460815.1 AW460815 da25c08.y1 Xenla 13LiC1 Xenopus laevis cD...	225	8e-59
gb BE992046.1 BE992046 UI-M-B21-bec-d-07-0-UI.s1 NIH_BMAP_MHI2 S...	225	8e-59
gb BI839734.1 BI839734 fq42d11.y1 zebrafish adult brain Danio re...	225	8e-59
gb BI429515.1 BI429515 fr70h03.y1 zebrafish adult brain Danio re...	225	8e-59
gb BI472934.1 BI472934 fr93f12.y1 zebrafish adult brain Danio re...	225	8e-59

Herkömmlicher Ansatz: Browsing

- Für jede verbleibende Sequenz
 - Komplette Proteinsequenz bei GenBank holen
- Wir haben oft geklickt
 - 1 + (Suche)
 - 14 + (Gendetails)
 - 14^*5 + (SwissProt)
 - 14^*5 + (Blast)
 - 14^*5^*X (Genbank)



Föderierter DBMS Ansatz

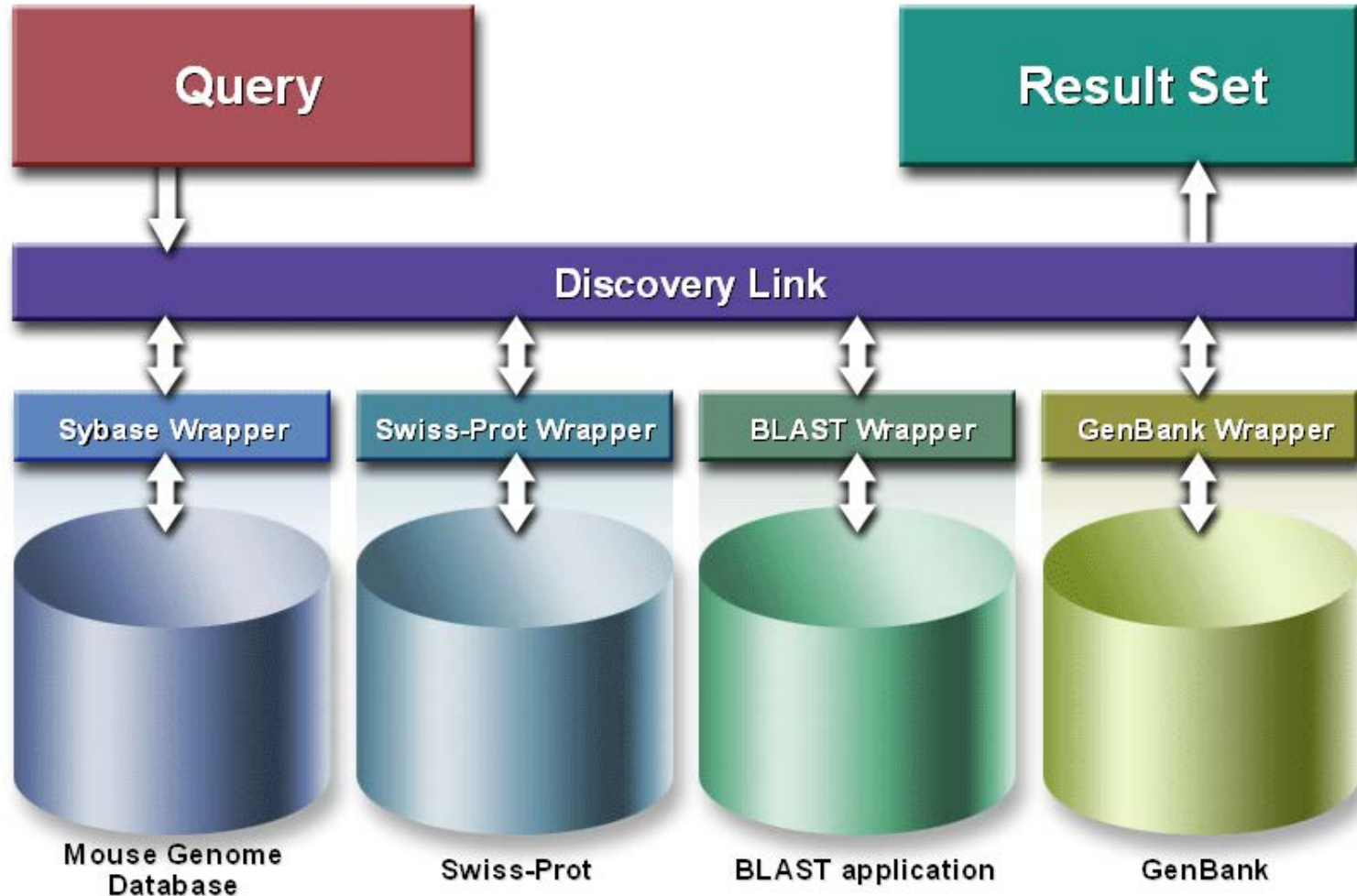
Finde alle menschlichen Sequenzen, die zu mindestens 60% identisch sind mit „channel“ Proteinen, die im Gewebe des zentralen Nervensystems der Maus exprimiert werden

- Ziel: Eine SQL-Anfrage

```
SELECT      g.accnum, g.sequence
FROM        genbank g, blast b, swissprot s, mgd m
WHERE       m.exp = "CNS"
AND         m.defn LIKE "%channel%"
AND         m.spid = s.id AND s.seq = b.query
AND         b.hit = g.accnum
AND         b.percentid > 60 AND b.species="HS";
```

- ohne alle Daten an einer Stelle zu speichern
 - On-the-fly integration

Beispiel: DiscoveryLink



Vergleich

DWH – Materialisierung	Föderation – Virtuelle Anfragen
Schnelle Anfragen (lokale Ausführung)	Langsame Anfragen (verteilte Ausführung)
Alle Anfragen	Nur beschränkte Anfragen (je nach API der Quellen)
Hoher Speicherbedarf	Kaum Speicherbedarf auf der Integrationsebene
Möglichkeit zur lokalen Änderung	Nur on-the-fly Korrekturen
Gefahr veralteter Daten, Divergenz zwischen Quelle und DWH	Stets aktuellste Daten
Daten müssen komplett verfügbar sein	Daten müssen über API zugreifbar sein
Prozedurale Integration	Deklarative Integration

Virtuelle Integration ist manchmal ein „Muss“

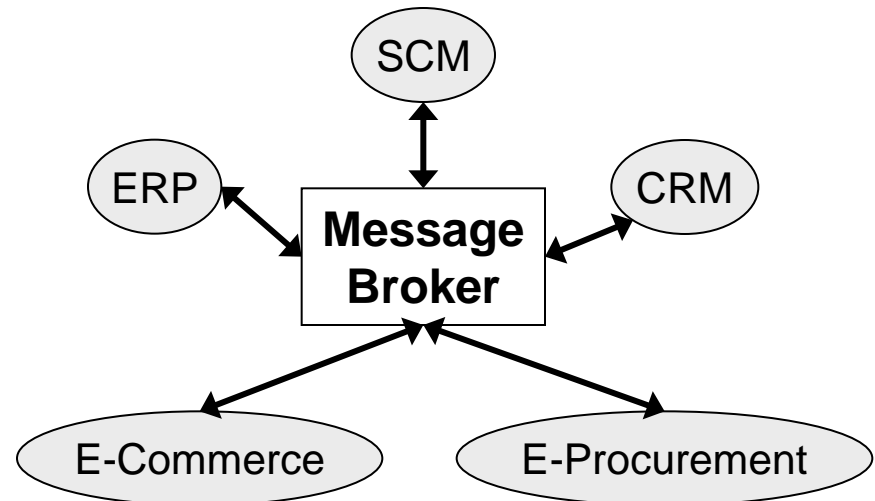
- Datenquellen zu groß (Web)
- Datenquellen nicht als Ganzes zugreifbar
 - Zugriffsbeschränkung, Copyrights, ...
- Inhalt der Datenquellen ändern sich sehr schnell
 - Börsenkurse, Newsticker, Preise, ...

Inhalt dieser Vorlesung

- Zwei Anwendungsszenarien
- **Abgrenzung**
 - Enterprise Application Integration
 - Objektorientierte Middleware
- Daten versus Schema

Enterprise Application Integration

- „**Integration ist ein Produkt**, kein Projekt“
- Viele kommerzielle Produkte und Anbieter
- Grundprinzip
 - Geschäftsvorfälle erzeugen Nachrichten
 - Diese werden an einen **Message Broker** gesendet
 - Der erkennt den Inhalt und wählt interessierte Quellen aus
 - Transformation der Nachrichten
 - Transaktionale Sicherheit („exactly once“)



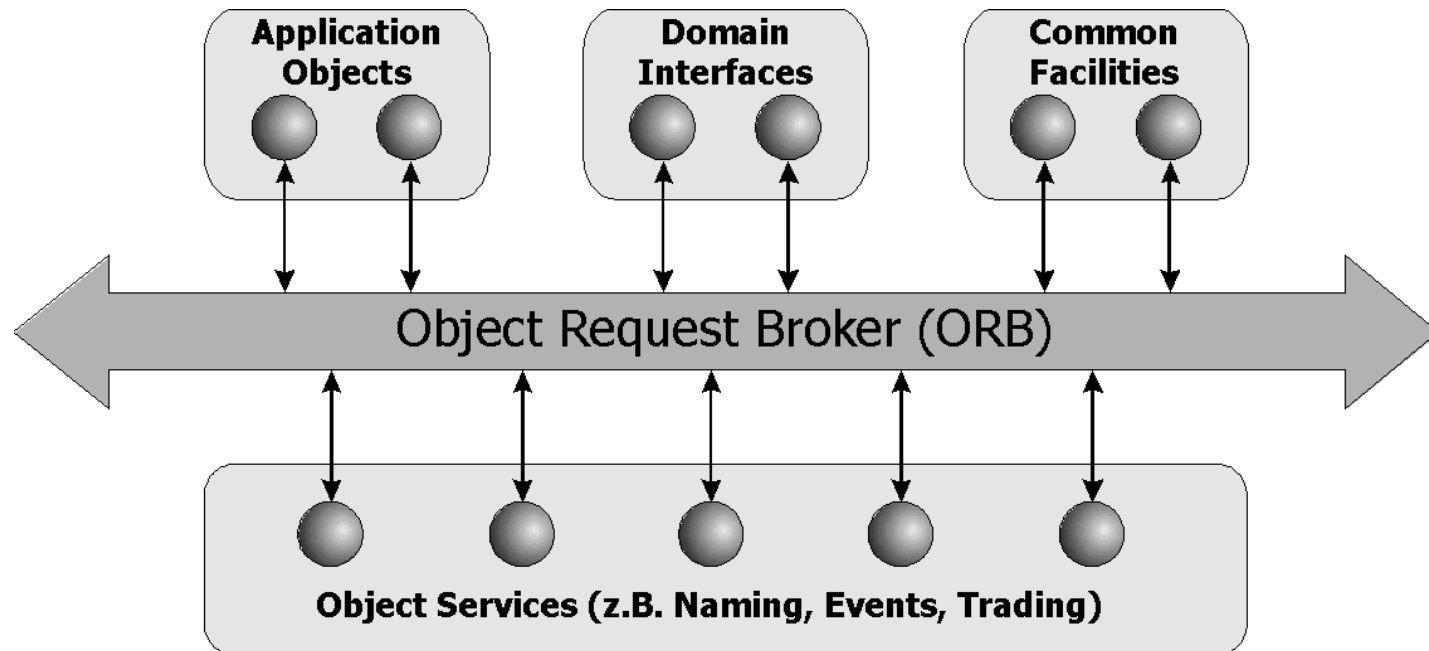
Unterschiede

- Enterprise Application Integration
 - Nachrichtenbasiert, keine Anfragen
 - Informationsverteilung
 - Aktion beim Eintreten eines Ereignisses
 - „Business Process Integration“
- Informationsintegration
 - Anfragebasiert
 - Annahme eines (praktisch) statischen Datenbestands
 - Aktion
 - Erst bei Anfrage (virtuelle Integration)
 - In regelmäßigen Zyklen (materialisierte Integration)
 - Datenbankintegration, Datenintegration

(OO) Middleware

- „Verteilung von Objekten wird vollkommen transparent“
- Viele Ansätze: DCOM, CORBA, OpenView, J2EE, ...
 - Weiterentwicklung von Remote Procedure Calls
- Grundidee
 - Objekte haben weltweit **eindeutige ID und Interfaces**
 - Clients programmieren gegen das Interface
 - Broker finden aufgerufene Objekte zur Laufzeit anhand ID
 - Marshalling, RPC, unmarshalling
 - Ziel: **Plattformunabhängigkeit**
 - Sprache, Betriebssystem, Kodierung, Protokoll

Beispiel: CORBA



- **Verteilungstransparenz**
- Domänenspezifische und generische **Dienste**
 - Namensauflösung, Nachrichtenübermittlung, Persistenz, ...

Unterschiede

- Middleware
 - **Werkzeug für die Programmerstellung**
 - Basiert auf Funktionsaufrufen
 - Versendung programmiersprachlicher Objekte
 - Vorwiegend transiente Objekte
 - Fokus auf **Methodenaufrufen**
- Informationsintegration
 - **Ist eine Aufgabe (Projekt)**, kein Werkzeug
 - Basiert auf Anfragen
 - Arbeiten mit strukturierten Daten
 - **Vorwiegend persistente Objekte**
 - Fokus auf Objektmanipulation und -suche

Synergie

- Middleware: Plattform zur Entwicklung von integrierten Systemen
 - Lösung der Verteilungsproblematik
 - Auflösung von low-level Formatunterschieden (ASCII versus UniCode, Little Endian versus Big Endian, Datumsformate, ...)
 - Auflösung von Unterschieden im Betriebssystem, Programmiersprache, ...
- Middleware leistet i.A. nicht
 - Umgang mit strukturierten Daten / Schemata
 - Übersetzung von Anfragen
 - **Semantische Integration**
- EAI wird i.d.R. auf Middleware implementiert

Inhalt dieser Vorlesung

- Zwei Anwendungsszenarien
- Abgrenzung
- Daten versus Schema
 - Extension versus Intension
 - Redundanz versus Komplementierung

Intension & Extension

- Intension
 - Die Intension eines Informationssystems ist die Menge der **Schemainformationen** und deren **Bedeutung (Semantik)**
 - Schemata, Metadaten, Modelle, Klassen
- Extension
 - Die Extension eines Informationssystems ist die Menge aller in ihm **gespeicherten Daten**
 - Instanzen, Tupel, Daten, Objekte, Entitäten
- Die Intension gibt der Extension Struktur
- Informationsintegration muss beides betrachten
 - Intension: Schemaheterogenität
 - Extension: Widersprüchliche Daten

Intension und Extension einer Tabelle

- **Intension**

- Struktur einer Menge von Entitäten
- Semantik der Struktureinheiten
- Statisch

- **Extension**

- Zustand der Tabelle
- Menge von Entitäten
- Dynamisch

Buch		
ISBN	Titel	Autor
3442727316	Moby Dick	Herman Melville
3491960827	Robinson Crusoe	Daniel Defoe
3462032283	Zwölf	Nick McDonell
3883891606	Timbuktu	Paul Auster
...

Redundanz und Komplementierung

- Redundanz (Überlappung)
 - In Extension und Intension möglich
 - Segen: Ohne **minimale Redundanz** ist Integration meist sinnlos
 - Was gehört zu was?
 - Fluch: Führt zu **Widersprüchen**, Doppelungen, ...
- Komplementierung
 - Informationen mehrerer Quellen werden zu einem **größeren Ganzen** integriert
 - Der eigentliche Sinn von Informationsintegration

Intensionale Redundanz

Quelle 1

ISBN	Author	Pages
3442727316	Herman Melville	1056
978-3491960824	Daniel Defoe	644
3462032283	Nick McDonell	240
3883891606	Paul Auster	227

Quelle 2

ISBN	Autorname	Year
3491960827	Daniel Defoe	1719
3442727316	H Melville	1851
3462026496	Saul Bellow	1992

Extensionale Redundanz

ISBN	Author	Pages
3442727316	Herman Melville	1056
978- 3491960824	Daniel Defoe	644

ISBN	Autorname	Year
3491960827	Daniel Defoe	1719
3442727316	H Melville	1851

- Extensionale Redundanz: Menge der von zwei Quellen **gemeinsam repräsentierten Objekte** ist nicht leer
- Voraussetzung dafür ist intensionale Redundanz
 - Gleiche Objekte müssen aus der gleichen „Klasse“ sein
- Oftmals schwer zu erkennen (Duplikaterkennung)

Probleme

ISBN	Author	Pages
3442727316	Herman Melville	1056
978-3491960824	Daniel Defoe	844

ISBN	Autorname	Year
3491960827	Daniel Defoe	1719
3442727316	H Melville	1851

Schwer zu
erkennendes
Duplikat

Datenkonflikt

Komplementierung (nach Datenfusion)

Intensionale Komplementierung

Extensionale
Komplementierung

ISBN	Author	Pages	Year
344272731 6	Herman Melville	1056	1851
349196082 4	Daniel Defoe	644	1719
346203228 3	Nick McDonell	240	
388389160 6	Paul Auster	227	
346202649 6	Saul Bellow		1992

Zusammenfassung

- Intensionale Redundanz ermöglicht extensionale Komplementierung
 - Zwei Quellen mit teilweise „gleichem“ Schema können zu einer **überdeckenderen** Datenbasis integriert werden
 - Coverage
- Extensionale Redundanz ermöglicht intensionale Komplementierung
 - Zwei Quellen, die über gleiche Objekte sprechen, können zu einer **dichteren** Datenbasis integriert werden
 - Density
- Insgesamt ist das Ziel der Integration eine **vollständigere** Datenbasis (completeness)