

Relational Databases for Biologists: Efficiently Managing and Manipulating Your Data

Session 1 Data Conceptualization and Database Design

Robert Latek, Ph.D.
Sr. Bioinformatics Scientist
Whitehead Institute for Biomedical Research

WIBR Bioinformatics, © Whitehead Institute, 2004

What is a Database?

- A collection of data
- A set of rules to manipulate data
- A method to mold information into knowledge
- Is a phonebook a database?
 - Is a phonebook with a human user a database?

| | | |
|-------------|---------------------------|----------|
| Babbitt, S. | 38 William St., Cambridge | 555-1212 |
| Baggins, F. | 109 Auburn Ct., Boston | 555-1234 |
| Bayford, A. | 1154 William St., Newton | 555-8934 |

WIBR Bioinformatics, © Whitehead Institute, 2004

Why are Databases Important?

- Data -> Information -> Knowledge
- Efficient Manipulation of Large Data Sets
- Integration of Multiple Data Sources
- Cross-Links/References to Other Resources

WIBR Bioinformatics, © Whitehead Institute, 2004

Why is a Database Useful?

- If Database Systems Simply Manipulate Data, Why not Use Existing File System and Spreadsheet Mechanisms?
- “Baggins” Telephone No. Lookup:
 - Human: Look for B, then A, then G ...
 - Unix: `grep Baggins boston_directory.txt`
 - DB: `SELECT * FROM directory WHERE IName="Baggins"`

| | | |
|-------------|---------------------------|----------|
| Babbitt, S. | 38 William St., Cambridge | 555-1212 |
| Baggins, F. | 109 Auburn Ct., Boston | 555-1234 |
| Bayford, A. | 1154 William St., Newton | 555-8934 |

WIBR Bioinformatics, © Whitehead Institute, 2004

What is the Advantage of a Database?

- Find All Last Names that Contain “Th” but do not have Street Address that Begin with “Th”.
 - Human: Good Luck!
 - UNIX: Write a directory parser and a filter.
 - DB: `SELECT IName FROM directory WHERE IName LIKE "%th%" AND street NOT LIKE "Th%"`

WIBR Bioinformatics, © Whitehead Institute, 2004

Why Biological Databases?

- Too Much Data
- Managing Experimental Results
- Improved Search Sensitivity
- Improved Search Efficiency
- Joining of Multiple Data Sets

WIBR Bioinformatics, © Whitehead Institute, 2004

Still Not Convinced?

- The Typical Excel Spreadsheet of Microarray Data

| Affy | lung | cardiac | gall_bladder | pancreas | testis |
|----------|------|---------|--------------|----------|--------|
| 92632_at | 20 | 20 | 20 | 20 | 20 |
| 94246_at | 20 | 71 | 122 | 20 | 20 |
| 93645_at | 216 | 249 | 152 | 179 | 226 |
| 98132_at | 135 | 236 | 157 | 143 | 145 |

- Now Find All of the Genes that have 2-3 fold Over-Expression in the Gall Bladder Compared to the Testis

WIBR Bioinformatics, © Whitehead Institute, 2004

Mini-Course Goals

- Conceptualize Data in Terms of Relations (Database Tables)
- Design Relational Databases
- Use SQL Commands to Extract/Data Mine Databases
- Use SQL Commands to Build and Modify Databases

WIBR Bioinformatics, © Whitehead Institute, 2004

Session Outline

- Session 1
 - Database background and design
- Session 2
 - SQL to data mine a database
- Session 3
 - SQL to create and modify a database
- Demonstration and Lab

WIBR Bioinformatics, © Whitehead Institute, 2004

Supplemental Information

- <http://jura.wi.mit.edu/bio/education/bioinfo-mini/db4bio/>
- <http://www.mysql.com/documentation/>
- A First Course In Database Systems. Ullman and Widom .
 - ISBN:0-13-861337-0

WIBR Bioinformatics, © Whitehead Institute, 2004

Flat vs. Relational Databases

- Flat File Databases Use Identity Tags or Delimited Formats to Describe Data and Categories Without Relating Data to Each Other
 - Most biological databases are flat files and require specific parsers and filters
- Relational Databases Store Data in Terms of Their Relationship to Each Other
 - A simple query language can extract information from any database

WIBR Bioinformatics, © Whitehead Institute, 2004

GenBank Report

```
LOCUS H2-K 1585 bp mRNA linear ROD 15-NOV-2002
DEFINITION Mus musculus histocompatibility 2, K region (H2-K), mRNA.
ACCESSION XM_193866
VERSION XM_193866.1 GI:29544196
KEYWORDS .
SOURCE Mus musculus (house mouse)
ORGANISM Mus musculus.
REFERENCE 1 (bases 1 to 1585)
AUTHORS NCBI Annotation Project.
TITLE Direct Submission
JOURNAL Submitted (13-NOV-2002) National Center for Biotechnology
COMMENT GENOME ANNOTATION RESEQ
FEATURES             Location/Qualifiers
     source            1..1585
                        /organism="Mus musculus"
                        /strain="C57BL/6J"
                        /db_xref="taxon:10090"
                        /chromosome="17"
     gene              1..1585
                        /gene="H2-K"
                        /db_xref="LocusID:14972"
                        /db_xref="MGI:95904"
     CDS                223..1137
                        /gene="H2-K"
                        /codon_start=1
                        /product="histocompatibility 2, K region"
                        /protein_id="XP_193866.1"
                        /translation="MSGRGRGOWNSRRPDSIGSRHRKPRMSRVSEWTLRT...
BASE COUNT  350 a 423 c 460 g 352 t
ORIGIN
1 gaagtcgqga atcgcagaca gttgcgtagg taccgtgac gctgctctg ctgtggcgg
WIBR Bioinformatics, © Whitehead Institute, 2004
```

NCBI NR Database File

```
>gij2137523|pir||59068 MHC class I H2-K-b-alpha-2 cell surface glycoprotein - mouse (fragment)
AHTIQVISGCEVGS DGRLLRGYQQYAYDGC DYALNEDLKTWTAADMAALITKHKWEQAGEAERL RAYLE
GTCVEWLRRLYKNGNATLLRT

>gij25054197|ref|XP_193866.1| histocompatibility 2, K region [Mus musculus]
MSRGRGGWSRRG P S I G S G R H R K P R A M S R V S E W T L R T L L G Y Y N Q S K G G S H T I Q V I S G C E V G S D G R L L R G Y
Q
QYAYDGC DYALNEDLKTWTAADMAALITKHKWEQAGEAERL RAYLE G T C V E W L R R Y L K N G N A T L L R T D S
PKAHVTHHSR P E D K V T L R C W A L G F Y P A D I L T W Q L N G E E L I Q D M E L V E T R P A G D G T F Q K W A S V V P L G K E
Q Y Y T C H V Y H Q G L P E P L T L R W E P P P S T V S N M A T V A V L V L G A A I V T G A V V A F V M K M R R R N T G G K G G D Y A L A
P G S Q T S D L S L P D C K V M V H D P H S L A

>gij25032382|ref|XP_207061.1| similar to histocompatibility 2, K region [Mus musculus]
MVPCTLLLLLAAALAPTQT R A G P H S L R Y F V T A V S R P L G E P R Y M E V G Y V D D E T F V R F D S A E N P R Y E P R A
R W M E Q E G P E Y W E R E T Q K A K G N E Q S F R V D L R T L L G Y Y N Q S K G G S H T I Q V I S G C E V G S D G R L L R G Y Q Q Y A Y
D
G C D Y I A L N E D L K T W T A A D M A A L I T K H K W E Q A G E A E R L R A Y L E G T C V E W L R R Y L K N G N A T L L R T D S P K A H V
T H H S R P E D K V T L R C W A L G F Y P A D I L T W Q L N G E E L I Q D M E L V E T R P A G D G T F Q K W A S V V P L G K E Q Y Y T C
H Y Y H Q G L P E P L T L R W E P P P S T V S N M A T V A V L V L G A A I V T G A V V A F V M K M R R R N T G G K G G D Y A L A P G S Q T
S D L S L P D C K V M V H D P H S L A
```

WIBR Bioinformatics, © Whitehead Institute, 2004

WIBR Bioinformatics, © Whitehead Institute, 2004

The Relational Database

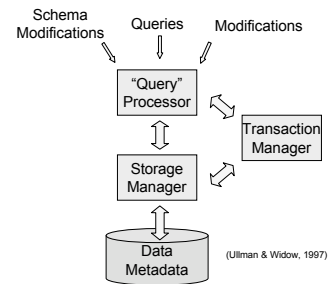
- Data is Composed of Sets of Tables and Links
- Structured Query Language (SQL) to Query the Database
- DBMS to Manage the Data

DBMS

- Database Management System (ACID)
 - Atomicity: Data independence
 - Consistency: Data integrity and security
 - Isolation: Multiple user accessibility
 - Durability: Recovery mechanisms for system failures

WIBR Bioinformatics, © Whitehead Institute, 2004

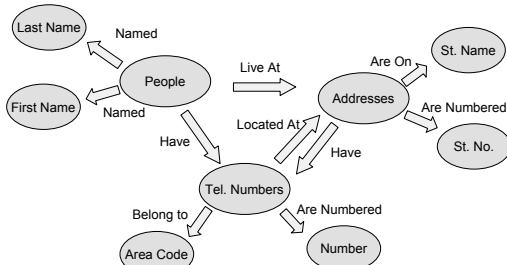
DBMS Architecture



WIBR Bioinformatics, © Whitehead Institute, 2004

Data Conceptualization

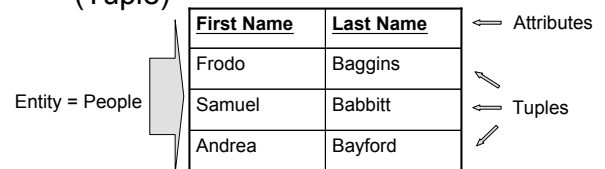
- Data and Links (For a Phonebook)



WIBR Bioinformatics, © Whitehead Institute, 2004

Data Structure

- Data Stored in Tables with Multiple Columns (Attributes).
- Each Record is Represented by a Row (Tuple)



WIBR Bioinformatics, © Whitehead Institute, 2004

Relational Database Specifics

- Tables are Relations
 - You perform operations on the tables
- No Two Tuples (Rows) can be Identical
- Each Attribute for a Tuple has only One Value
- Tuples within a Table are Unordered
- Each Tuple is Uniquely Identified by a Primary Key

WIBR Bioinformatics, © Whitehead Institute, 2004

Primary Keys

- Primary Identifiers (Ids)
- Set of Attributes that Uniquely Define a Single, Specific Tuple (Record)
- Must be Absolutely Unique
 - SSN ?
 - Phone Number ?
 - ISBN ?

| First Name | Last Name | SSN |
|------------|-----------|-------------|
| Frodo | Baggins | 332-97-0123 |
| Frodo | Binks | 398-76-5327 |
| Maro | Baggins | 215-01-3965 |

WIBR Bioinformatics, © Whitehead Institute, 2004

Find the Keys

| First Name | Last Name | SSN | Phone Number | Address |
|------------|-------------|-------------|--------------|----------------|
| Frodo | Baggins | 321-45-7891 | 123-4567 | 29 Hobbitville |
| Aragon | Elf-Wantabe | 215-87-7458 | 258-6109 | 105 Imladris |
| Boromir | Ringer | 105-91-0124 | 424-9706 | 31 Hobbitville |
| Bilbo | Baggins | 198-02-2144 | 424-9706 | 29 Hobbitville |
| Legolas | Elf | 330-78-4230 | 555-1234 | 135 Imladris |

WIBR Bioinformatics, © Whitehead Institute, 2004

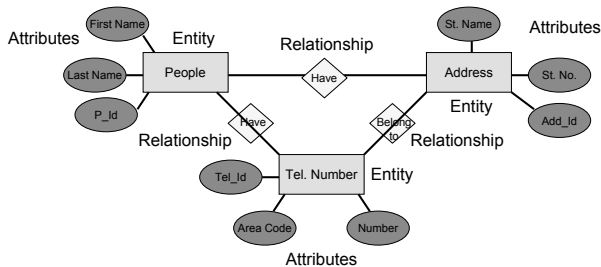
Design Principles

- Conceptualize the Data Elements (Entities)
- Identify How the Data is Related
- Make it Simple
- Avoid Redundancy
- Make Sure the Design Accurately Describes the Data!

WIBR Bioinformatics, © Whitehead Institute, 2004

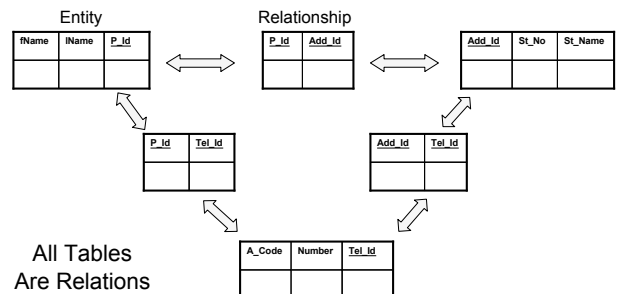
Entity-Relationship Diagrams

- Expression of a Database Table Design



WIBR Bioinformatics, © Whitehead Institute, 2004

E-R to Table Conversion



WIBR Bioinformatics, © Whitehead Institute, 2004

Steps to Build an E-R Diagram

- Identify Data Attributes
- Conceptualize Entities by Grouping Related Attributes
- Identify Relationships/Links
- Draw Preliminary E-R Diagram
- Add Cardinalities and References

WIBR Bioinformatics, © Whitehead Institute, 2004

Developing an E-R Diagram

- Convert a GenBank File into an E-R Diagram

```

LOCUS IL2RG 1451 bp mRNA linear PRI 17-JAN-2003
DEFINITION Homo sapiens interleukin 2 receptor, gamma (severe combined immunodeficiency) (IL2RG), mRNA.
ACCESSION NM_000206
VERSION NM_000206.1 GI:4557881
ORGANISM Homo sapiens
REFERENCE 1 (bases 1 to 1451)
AUTHORS Takeshita,T, Asano,H, Ohtani,K, Ishii,N, Kumaki,S, Tanaka,N, Munakata,H, Nakamura,M and Sugamura,K
TITLE Cloning of the gamma chain of the human IL-2 receptor
JOURNAL Science 257 (5068), 379-382 (1992)
MEDLINE 9233983
PUBMED 1631559
REFERENCE 2 (bases 1 to 1451)
AUTHORS Noguichi,M, YJH, Rosenblatt,H.M., Filipovich,A.H., Adelstein,S., Modi,W.S., McBride,O.W. and Leonard,W.J.
TITLE Interleukin-2 receptor gamma chain mutation results in X-linked severe combined immunodeficiency in humans
JOURNAL Cell 73 (1), 147-157 (1993)
MEDLINE 93214986
PUBMED 8462096
CDS 15..1124
/gene="IL2RG"
/product="interleukin 2 receptor, gamma chain, precursor"
/protein_id="WP_000197.1"
/db_xref="GI:4557882"
/db_xref="LocusID:3561"
/translation="MLKPSLPFTSLFLQLPLLVGLNLTILTPNGNEDTADFFLT..."
BASE COUNT 347 a 422 c 313 g 369 t
ORIGIN
1 gaagagcaag cgccatgttg aagccatcat taccattcac atccctctta ttccgcagc
    
```

WIBR Bioinformatics, © Whitehead Institute, 2004

Identify Attributes

- Locus, Definition, Accession, Version, Source Organism
- Authors, Title, Journal, Medline Id, PubMed Id
- Protein Name, Protein Description, Protein Id, Protein Translation, Locus Id, GI
- A count, C count, G count, T count, Sequence

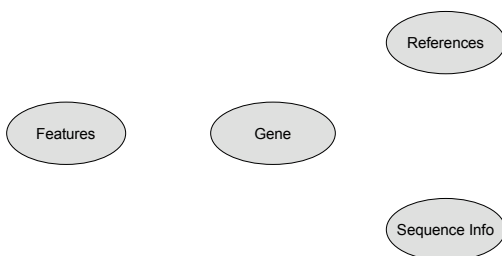
WIBR Bioinformatics, © Whitehead Institute, 2004

Identify Entities by Grouping

- Gene
 - Locus, Definition, Accession, Version, Source Organism
- References
 - Authors, Title, Journal, Medline Id, PubMed Id
- Features
 - Protein Name, Protein Description, Protein Id, Protein Translation, Locus Id, GI
- Sequence Information
 - A count, C count, G count, T count, Sequence

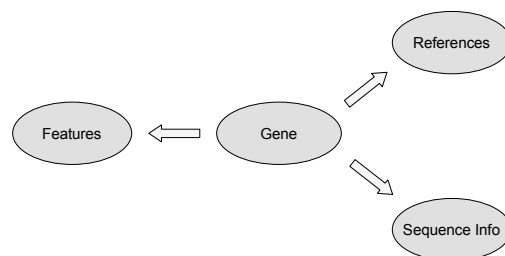
WIBR Bioinformatics, © Whitehead Institute, 2004

Conceptualize Entities



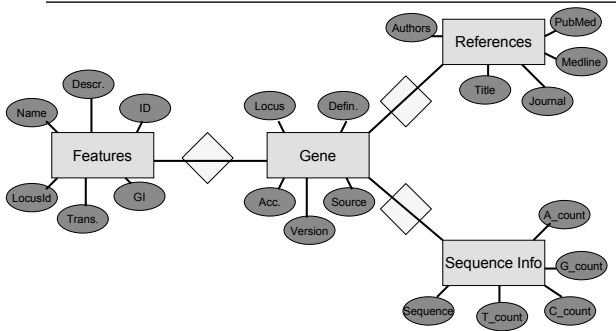
WIBR Bioinformatics, © Whitehead Institute, 2004

Identify Relationships



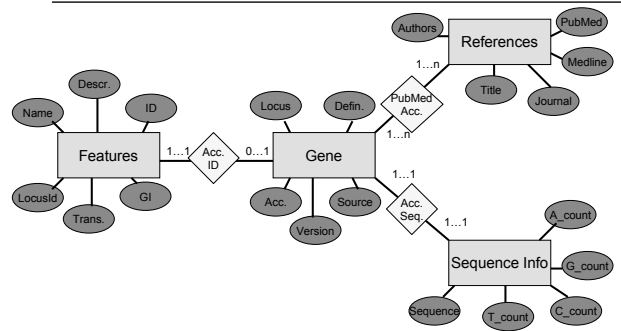
WIBR Bioinformatics, © Whitehead Institute, 2004

Preliminary E-R Diagram



WIBR Bioinformatics, © Whitehead Institute, 2004

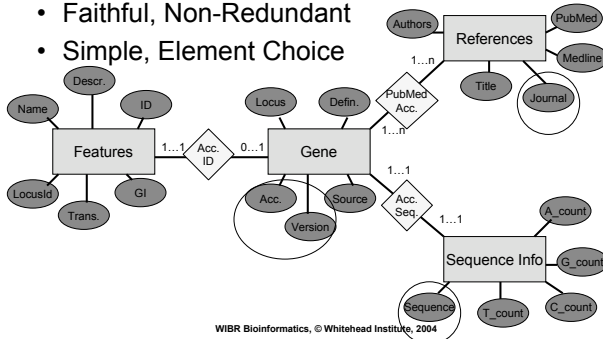
Cardinalities and References



WIBR Bioinformatics, © Whitehead Institute, 2004

Apply Design Principles

- Faithful, Non-Redundant
- Simple, Element Choice



WIBR Bioinformatics, © Whitehead Institute, 2004

Build Your Own E-R Diagram

- Express the Following Annotated Microarray Data Set as an E-R Diagram

| Affyid | GenBankId | Name | Description | LocusLinkId | LocusDescr | NT | RefSeq | AA | RefSeq |
|--------------|-----------|--------|-------------|-------------|------------|-----------|-----------|----|--------|
| U95-32123_at | L02870 | COL7A1 | Collagen | 1294 | Collagen | NM_000094 | NP_000085 | | |
| U98-40474_at | S75295 | GBE1 | Glucan | 2632 | Glucan | NM_000158 | NP_000149 | | |

| UnigeneId | GO Acc. | GO Descr. | Species | Source | Level | Experiment |
|-----------|---------|--------------|---------|----------|-------|------------|
| Hs.1640 | 0005202 | Serine Prot. | Hs | Pancreas | 128 | 1 |
| Hs.1691 | 0003844 | Glucan Enz. | Hs | Liver | 57 | 2 |

WIBR Bioinformatics, © Whitehead Institute, 2004

Summary

- Databases Provide ACID
- Databases are Composed of Tables (Relations)
- Relations are Entities that have Attributes and Tuples
- Databases can be Designed from E-R Diagrams that are Easily Converted to Tables
- Primary Keys Uniquely Identify Individual Tuples and Represent Links between Tables

WIBR Bioinformatics, © Whitehead Institute, 2004

Next Week

- Using Structured Query Language (SQL) to Data Mine Databases
- SELECT a FROM b WHERE c = d
- 5th Floor Conference Room on Monday, February 10.

WIBR Bioinformatics, © Whitehead Institute, 2004

Identify Attributes

| AffyId | GenBankId | Name | Description | LocusLinkId | LocusDescr | NT_RefSeq | AA_RefSeq | \\ |
|--------------|-----------|--------|-------------|-------------|------------|-----------|-----------|----|
| U95-32123_at | L02870 | COL7A1 | Collagen | 1294 | Collagen | NM_000094 | NP_000085 | \\ |
| U98-40474_at | S75295 | GBE1 | Glucan | 2632 | Glucan | NM_000158 | NP_000149 | \\ |

| UnigenId | GO Acc. | GO Descr. | Species | Source | Level | Experiment |
|----------|---------|--------------|---------|----------|-------|------------|
| Hs.1640 | 0005202 | Serine Prot. | Hs | Pancreas | 128 | 1 |
| Hs.1691 | 0003844 | Glucan Enz. | Hs | Liver | 57 | 2 |

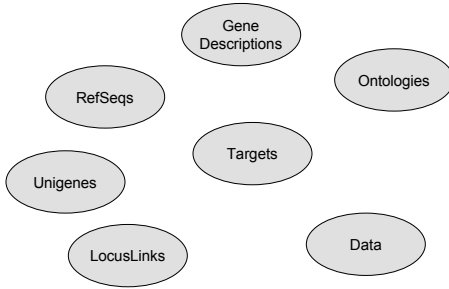
WIBR Bioinformatics, © Whitehead Institute, 2004

Identify Entities by Grouping

- Gene Descriptions
 - Name, Description, GenBank
- RefSeqs
 - NT RefSeq, AA RefSeq
- Ontologies
 - GO Accession, GO Terms
- LocusLinks
- Unigenes
- Data
 - Sample Source, Level
- Targets
 - Affy ID, Experiment Number, Species

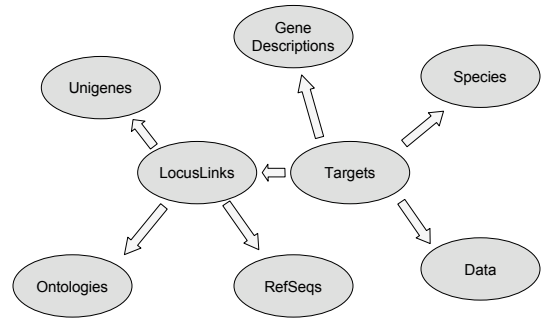
WIBR Bioinformatics, © Whitehead Institute, 2004

Conceptualize Entities



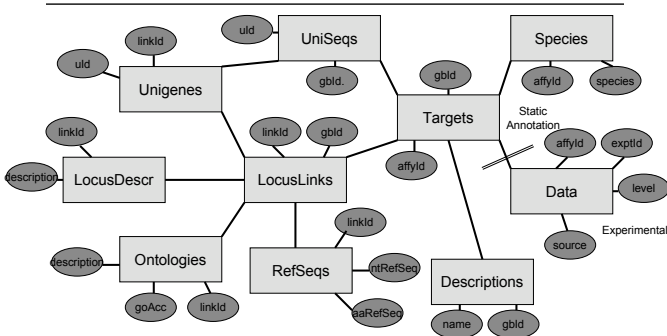
WIBR Bioinformatics, © Whitehead Institute, 2004

Identify Relationships



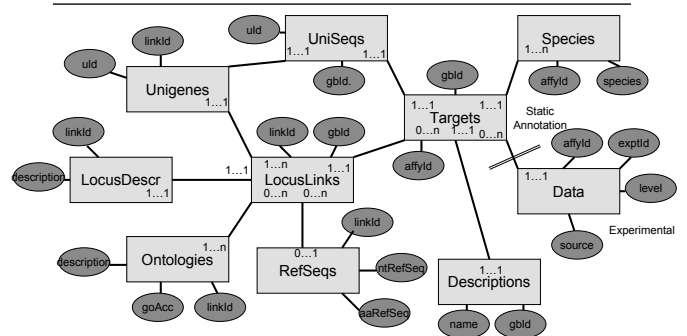
WIBR Bioinformatics, © Whitehead Institute, 2004

Preliminary E-R Diagram



WIBR Bioinformatics, © Whitehead Institute, 2004

Cardinalities and References



WIBR Bioinformatics, © Whitehead Institute, 2004

Apply Design Principles

