

Data Compression & Databases

Practical Introduction to Programming for
Scientists - Lecture 11

DNA Sequence

- seq="atggcagctaaagacgtaaaattcggtaacgacgctcgtgtgaaaatgctgcgcgccgtaaacgtactggcagatgca
gtgaaagttaccctcgggtccgaaaggccgtaacgtagttctggataaatcttcgggtgcaccgaccatcaccaagatgggtgttcc
gttgctcgtgaaatcgaactggaagacaagttcgaaaacatgggtgctgcagatgggtgaaagaagttgcctctaaagcgaacga
cgctgcaggcgacggtaccaccactgcaaccgtactggctcaggctatcatcactgaaggctgaaagctggtgctgcgggcat
gaacccgatggacctgaaacgtggtatcgacaaagctgttaccgctgcagttgaagaactgaaagcgtgtccgtaccgtgctct
gactctaaagcgattgctcaggttggtactatctccgctaactccgacgaaaccgtaggtaaactgatcgtgaagcgatggaca
aagtcggtaaagaaggcgttatcaccgttgaaagacggtaccggtctgcaggacgaactggacgtggttgaaaggtatgcagttcg
accgtggctacctgtctccttacttcatcaacaagccggaaactggcgcagtagaactggaaagcccgttcatcctgctggctga
caagaaaatctccaacatccgcgaaatgctgccggttctggaagccgttgccaaagcaggcaaacccgctgctgatcatcgctg
aagatgtagaaggcgaagcgtggcaactctggttgtaaacaccatgcgtggcatcgtgaaagttgctgcagttaaagctccggg
cttcggcgatcgtcgtaaagctatgctgcaggatatcgcaaccctgactggcgggtaccgtaatctctgaagagatcgggtatggag
ctggaaaagcaaccctggaagacctgggtcaggctaaacgcgttgatcaacaagacaccaccaccatcatcgatggcg
tgggcgaagaagctgcaatccagggccggtgtgctcagatccgtcagcagattgaagaagcaacttctgactacgaccgtgaa
aaactgcaggagcgcgtagcgaactggcaggcggcgttgacgttatcaaagtaggtgctgctaccgaagttgaaatgaaaga
gaaaaagcacgcgttgaaagacgccctgcacgcgaccctgctgcggtagaagaaggcgtggttgctgggtggtggtggtgctg
ctgatccgcgtagcgtctaaactggctgacctgctggtcagaacgaagaccagaacgtgggtatcaaagttgactgctgca
atggaagctccgctgcgtcagatcgtcctgaactgcggcgaagaaccgtctgttggtgctaacaccgttaaaggcggcgacggc
aactacggttacaacgcagcaaccgaagaatacggcaacatgatcgacatgggtatcctggaccaaccaaaagtaacccttc
tgctctgcagtacgcggcttctgtggctggcctgatgatcaccaccgaatgcatggttaccgacctgccgaaaaacgatgcagct
gacttaggcgctgctggcgggtatgggcggcatgggtggcatgggcggcatgatgtaa"

Data Compression

- gzip, zip, bzip2, ...
- `from zlib import compress, decompress`
- `compress(<str>, <level>)`
- `decompress(<cmpstr>)`
- great, but...

Compress DNA Seq

- With gzip, etc.
 - Slooow, 2.8 GB gator genome takes 10 min !
 - No random access of compressed data
- Can we come up with something better ?

Compress DNA Seq

- DNA, only 4 possible values -> 2 bits
- 1 Byte = 8 bits -> 4x compression !!!
- But...
- What about unknowns ?
- $5*5*5 = 125$
 - 3x compression
 - naturally ordered into triplets

```

import os

os.chdir("/tmp")

letmap={"a":1,"c":2,"g":3,"t":4,"n":0}

# build compression/decompression dict
tripmap={}
tripmapinv={}
for a in letmap.keys():
    for b in letmap.keys():
        for c in letmap.keys():
            val=chr(letmap[a]*25+letmap[b]*5+letmap[c])
            tripmap[a+b+c]=val
            tripmapinv[val]=a+b+c

print tripmap

seq=file("gator_small.seq","r")
outseq=file("gator_small.cmp","w")

while True:
    trip=seq.read(3)
    if len(trip)==0 : break
    try: outseq.write(tripmap[trip])
    except: break

```

```

#include <stdio.h>

int main(int argc, char *argv[]) {
    FILE *in,*out;

    in=fopen("/tmp/gator_small.seq","r");
    out=fopen("/tmp/gator_small.cmpr","w");

    char seq[3];
    char map[256];
    map['a']=1;
    map['c']=2;
    map['g']=3;
    map['t']=4;
    map['n']=0;
    while (fread(seq,3,1,in)==1) {
        unsigned char xlate=map[seq[0]]*25+map[seq[1]]*5+map[seq[2]];
        fwrite(&xlate,1,1,out);
    }
    fclose(in);
    fclose(out);
}

```



How to Store Complex Data ? (again)

- Students
 - Name, address, ...
- Classes
 - Description, Instructor, when offered, ...
- Class Year
 - Which class, year offered, students
- Grades
 - Class, student, grade

Spreadsheet !

Database Terms

- Table - equivalent to a single spreadsheet page
- Database - a set of tables
- Column - one specific type of data for all records in a table
- Row - one record in a table
- Tuple - another name for one row in the database
- Schema - The titles and datatypes of all the rows in all of the tables
- RDMS - Relational Database Management System

Relational Database

Tables (Relations)

First Name	Last Name	Program	Grade
John	Barnes	1	A
Frank	Smith	1	B+
Carol	Franklin	2	A-
Steve	Black	3	C+
Charles	Baker	1	B+

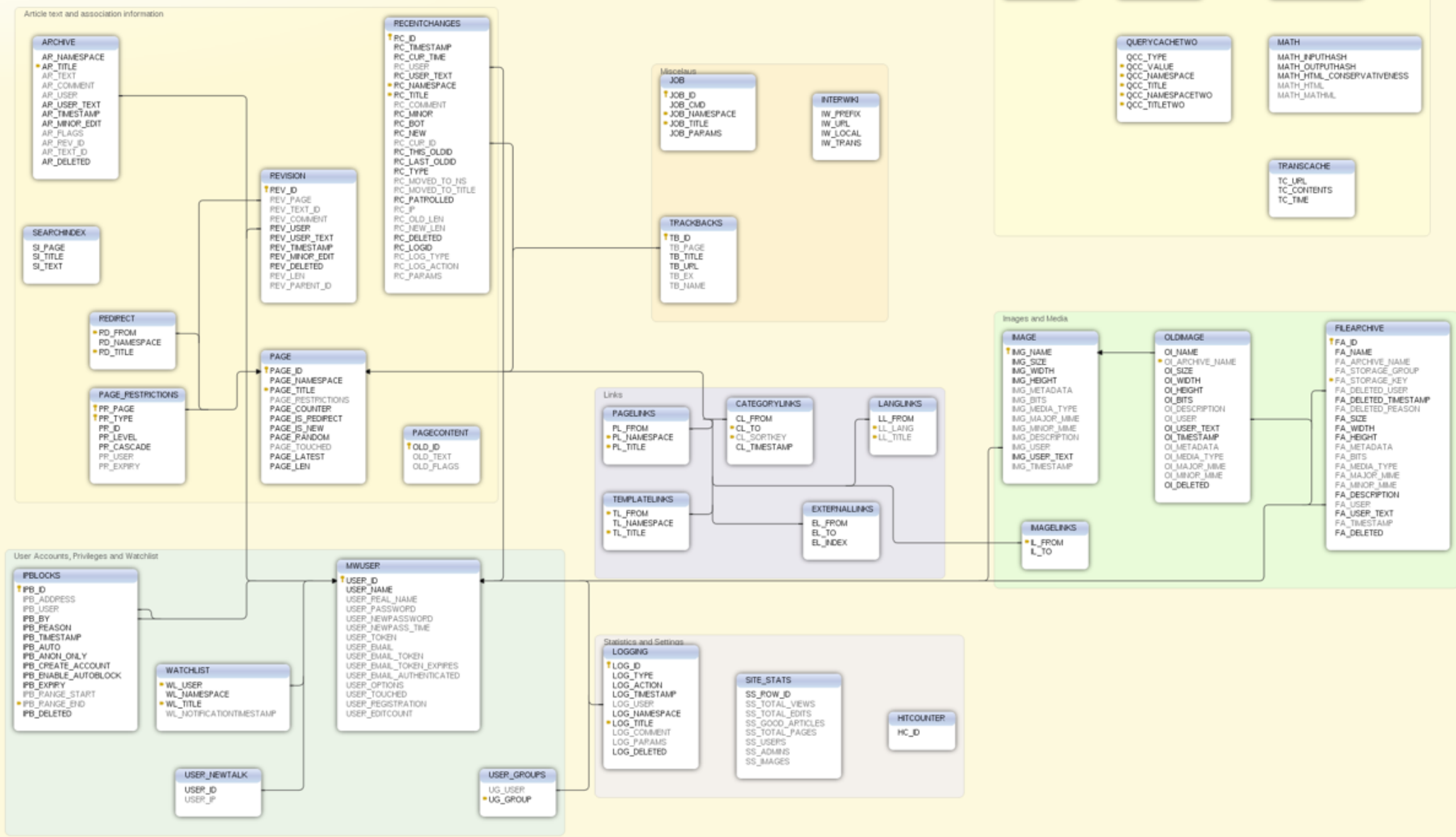
For many->many we use a 'junction table'.

1 to 1

id	Name	Institution
1	SCBMB	BCM
2	Biochemistry	BCM
3	Biochemistry	Rice

Database Schema

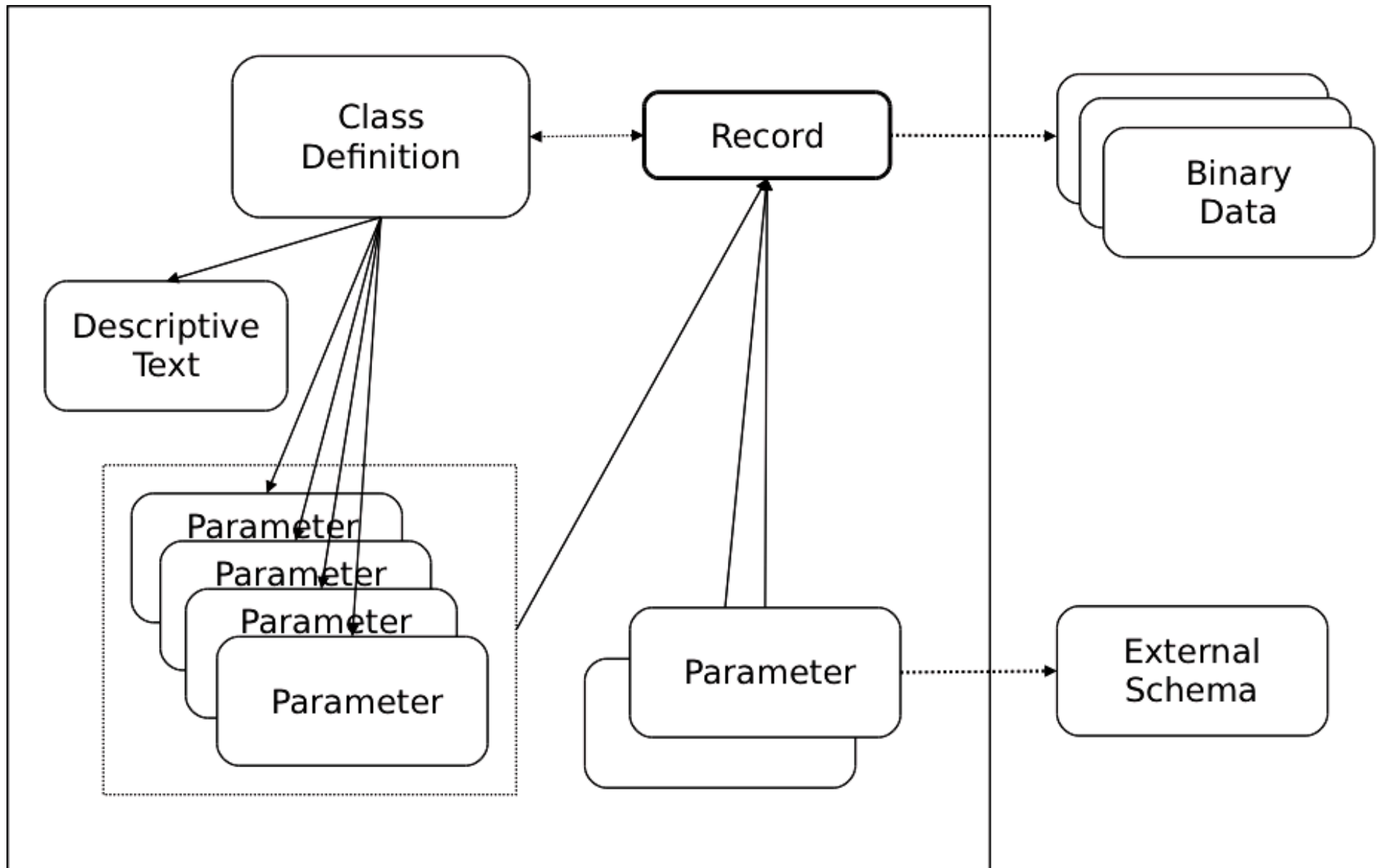
Media Wiki database schema
 Move the mouse cursor over the table header and columns to view the comments
 generated using DbSchema



Pythonic Approach

- ⑤ Dictionary
- ⑤ Class
- ⑤ List

OODB



Databases

- ⑤ Simple (embedded) Database
 - ⑤ dbm, BerkeleyDB (bsddb), SQLite*
- ⑤ Relational Database
 - ⑤ MySQL, Oracle, DB2, SQLite*
- ⑤ Object Oriented Database (OODB)
 - ⑤ Zope, Databeans
- ⑤ XML Database
 - ⑤ Sedna, Oracle, BerkeleyDB-XML

Simple/Embedded Database

- ④ Key/Value pairs
- ④ Python
 - ④ Shelve – Persistent dictionary
 - ④ anydbm – Generic access to DBM-style databases
 - ④ whichdb – Guess which DBM module created a database
 - ④ dbm – Simple “database” interface
 - ④ gdbm – GNU’s reinterpretation of dbm
 - ④ dbhash – DBM-style interface to the BSD database library
 - ④ dumbdbm – Portable DBM implementation
 - bsddb (bsddb3) – Interface to Berkeley DB library

pickle

- 'Serialization' - converting a complex object to a stream of data

```
from cPickle import dump,load,dumps,loads
```

```
dump(obj,file[,protocol]) # stores 'obj' in file
```

```
obj=load(file) # restores 'obj' from file
```

```
str=dumps(obj[,protocol]) # pickled representation of obj
```

```
obj=loads(str) # restore representation of obj
```

shelve

```
import shelve      # dictionary-like object on disk
dct=shelve.open(filename[,protocol])
dct=shelve.open(filename,writeback=True)
dct.close()
```

ACID

⌚ Atomicity

- ⌚ Transactions are 'all or none', no partial modifications

⌚ Consistency

- ⌚ Contents of the database are always valid and self-consistent

⌚ Isolation

- ⌚ Readers will not see partially completed transactions of writers

⌚ Durability

- ⌚ Completed transactions can survive any possible system crash

SQL

- ⑤ Structured Query Language
- ⑤ Developed at IBM in the 70s
- ⑤ ANSI Standard in 1986
- ⑤ Now fairly ubiquitous for Relational Databases
- ⑤ (see Wikipedia for more trivia)

SQLite3

```
import sqlite3
```

```
conn=sqlite3.connect(path)           # Create a database 'connection'
```

```
cursor=conn.cursor()                # Create a cursor object
```

```
cursor.execute(SQLstring)           # Execute a query
```

```
conn.commit()                        # After making any changes
```

```
cursor.close()                       # close the cursor when you're done
```

SQLite3

Python type	SQLite type
<u>None</u>	NULL
<u>int</u>	INTEGER
<u>long</u>	INTEGER
<u>float</u>	REAL
<u>str</u> (UTF8-encoded)	TEXT
<u>unicode</u>	TEXT
<u>buffer</u>	BLOB

SQLite3

- ④ `conn=sqlite3.connect("mydb.db")`
- ④ `cur=conn.cursor()`
- ④ create table
 - ④ `cur.execute("CREATE TABLE Person(Id INT, Name TEXT, Zip INT);")`
- ④ insert
 - ④ `cur.execute("INSERT INTO Person VALUES (1,'Steve',77884)")`
- ④ select
 - ④ `cur.execute("SELECT * FROM Person")`
 - ④ `print cur.fetchall()`
- ④ update
 - ④ `cur.execute("UPDATE Person SET zip=77584 WHERE Name='Steve'")`
- ④ drop
 - ④ `cur.execute("DROP TABLE Person")`

Presentations Monday!

- Remember to send me your class projects by midnight SATURDAY
- Follow all directions on the class website!
- Remember you will have only 5 minutes to present. Practice your talk in advance, and if possible, test the projector with your laptop.