BIOINFORMATICS
Perl Workshop

TWO PROBLEMS – PART II

GENOME
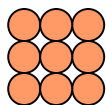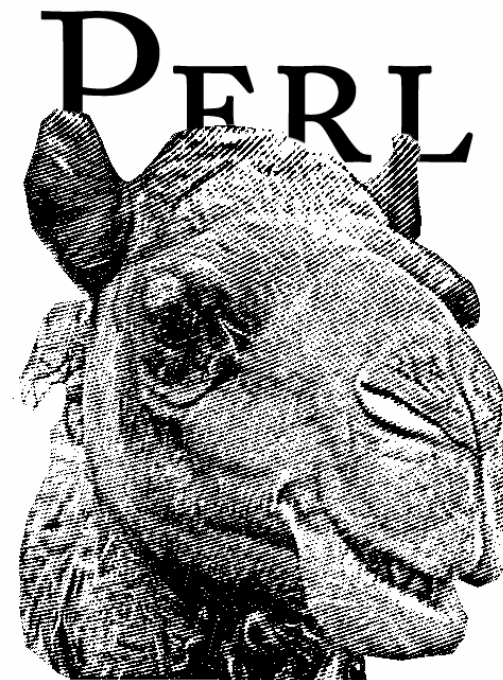SCIENCES
CENTRE

# Perl Panacea?

· The <u>camel</u> represents the desirable features of Perl
  · O'Reilly colophon

· Why is the camel successful?
  · adapted itself to desert environment
    · low water needs (gets around with what's around)
    · elegant from a distance
    · still comfortable
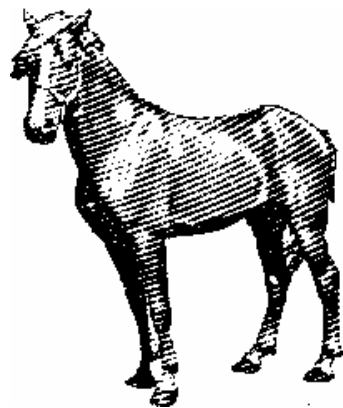  · not cute - until you get to know the camel

PERL

ADAPTIVE, LOW-MAINTENANCE
EASY TO RIDE

DATA OASIS                    APPLICATION OASIS

DESERT

BIOINFORMATICS
Perl Workshop

TWO PROBLEMS – PART II

GENOME
SCIENCES
CENTRE

# Perl as Explorer

Lots of camel mechanics in the desert, and we're in a desert

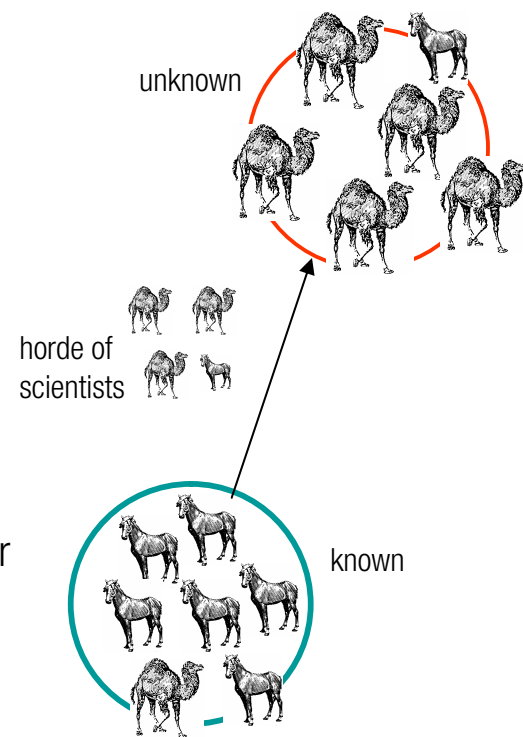| ANOTHER LEADING LANGUAGE | PERL | SIMPLIFIED EXPLORATION MODEL |
|---|---|---|



~ mathematics
beautiful & elegant
rigorous, requires overhead
fast on smooth ground
slow in rough terrain
! killed by camel veterinarian

~ physics
gets you there
explores uninhabited terrain, cavalier
average speed on smooth ground
average speed in rough terrain
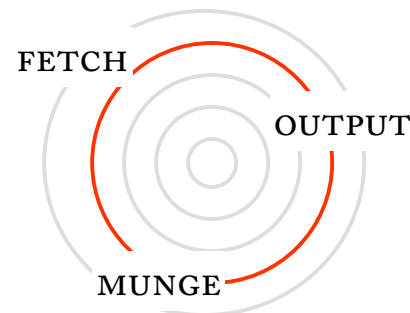* horse veterinarian OK

unknown

horde of scientists

known

BIOINFORMATICS
Perl Workshop

TWO PROBLEMS – PART II

GENOME
SCIENCES
CENTRE

# Holy Triad of Analysis

· many types of analyses fall into this <u>analysis triad</u>
  · fetch from: file, user, pipe, http, ftp
  · munge: collate, sort, organize, count, enumerate
  · output: text, image, HTML, XML
· each step is made <u>pleasant and easy</u> with Perl
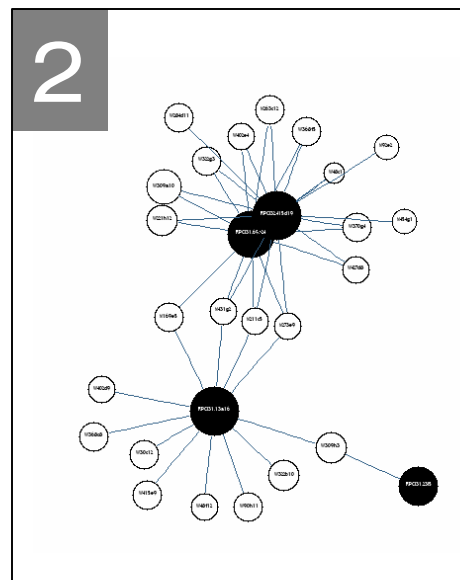
FETCH

OUTPUT

MUNGE

**1**

al Mapping Data, Nov/01/2002

DC-RatDataSet.txt  (TAB separated values, including RH-vectors, 1.8 MB)

No      - consecutive number
Chr     - chromosome
Cont    - contig number on respective chromosome
Lib.Clone - library name.clone name
F       - framework marker, according to Rat Genome Database, RGD
F cR    - cR-position of framework marker
YAC's   - YACs hit by Lib.Clone, "," separated
          MPMGy916 = ICRF/BWH/Liege_SHRSP Rat YAC Library
          WIBRY933 = Whitehead Institute / MIT rat YAC Library

Go to Chromosome: [1][2][3][4][5][6][7][8][9][10][11][12][13][14][15][16][17][18][19][20][21=X]

| No | Chr | Cont | Lib.Clone | F | F cR | |
|----|-----|------|-----------|---|------|---|
| 748 | 02 | 1 | RPCI31.64I13 | | | MPMGy916.186d9, MPMGy916.34f11, MPMGy916.528s8 |
| 749 | 02 | 1 | RPCI31.63c10 | | | MPMGy916.186d9, MPMGy916.345h7, MPMGy916.34f11, MPMGy916.528s8 |
| 750 | 02 | 1 | RPCI32.432g2 | | | MPMGy916.113f2, MPMGy916.16e3, MPMGy916.191c6, MPMGy916.195f6, MPMGy916.278b6, |
| 751 | | | | d2rat1 | 0.0 | |
| 752 | 02 | 1 | RPCI32.406p4 | | | MPMGy916.113f2, MPMGy916.16e3, MPMGy916.191c6, MPMGy916.195f6, MPMGy916.278b6, |
| 753 | 02 | 1 | RPCI31.49d41 | | | MPMGy916.113g2, MPMGy916.146d8, MPMGy916.172b7, MPMGy916.204a1, MPMGy916.271c |
| 754 | 02 | 1 | RPCI32.423c7 | | | MPMGy916.112f1, MPMGy916.113f2, MPMGy916.146d8, MPMGy916.172b7, MPMGy916.186g |
| 755 | 02 | 1 | RPCI31.168g3 | | | MPMGy916.112f1, MPMGy916.113g2, MPMGy916.146d8, MPMGy916.172b7, MPMGy916.204a |
| 756 | 02 | 1 | RPCI32.368g3 | | | MPMGy916.113g2, MPMGy916.135d3, MPMGy916.172b7, MPMGy916.219g3, MPMGy916.265 |
| 757 | 02 | 1 | RPCI32.420b8 | | | MPMGy916.113g2, MPMGy916.135d3, MPMGy916.172b7, MPMGy916.219g3, MPMGy916.265 |
| 758 | 02 | 1 | RPCI32.420c8 | | | MPMGy916.113g2, MPMGy916.135d3, MPMGy916.172b7, MPMGy916.219g3, MPMGy916.265 |
| 759 | 02 | 1 | RPCI31.69I13 | | | MPMGy916.113g2, MPMGy916.172b7, MPMGy916.219g3, MPMGy916.265b11, MPMGy916.306 |
| 760 | 02 | 1 | RPCI32.411o23 | | | MPMGy916.113g2, MPMGy916.172b7, MPMGy916.219g3, MPMGy916.265b11, MPMGy916.306 |
| 761 | 02 | 1 | RPCI31.30n10 | | | MPMGy916.113e2, MPMGy916.135d3, MPMGy916.172b7, MPMGy916.196b3, MPMGy916.256c |

**2**

**3**

```
0 0 RPCI31.80h3
0 1 MPMGy916.380h8
0 2 WIBRy933.259d6
...
0 45 WIBRy933.284b4
0 46 WIBRy933.219c12
0 47 MPMGy916.110d1
1 0 MPMGy916.369g12
1 1 MPMGy916.282g2
1 2 RPCI31.33m10
...
21 7 RPCI31.17n14
21 8 WIBRy933.106h8
21 9 RPCI31.17i17
```

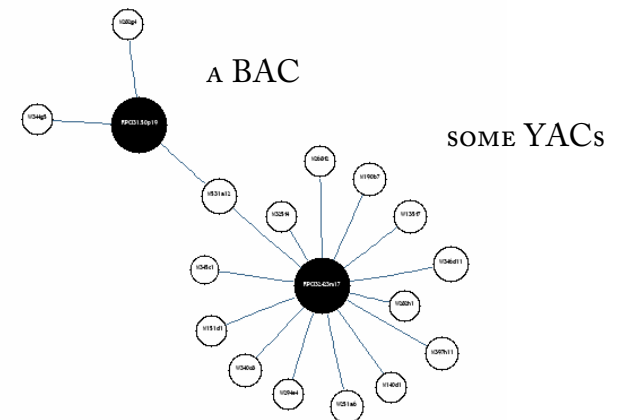! data we want is in a web table (<u>Very Bad Thing™</u>)

BAD
USAGE

* visualize the relationships for sanity

* format data to STDOUT

# Step 1 – Fetch – Perl Makes it Fun

| No | Clu | Cont | Lib.Clone | F | F cR | |
|----|-----|------|-----------|---|------|---|
| 748 | 02 | 1 | RPCI31.64l18 | | | MPMGy916.186d9, MPMGy916.34f11, MPMGy916.528b8 |
| 749 | 02 | 1 | RPCI31.63c10 | | | MPMGy916.186d9, MPMGy916.345h7, MPMGy916.34f11, MPMGy916.528b8 |
| 750 | 02 | 1 | RPCI32.432g2 | | | MPMGy916.113f2, MPMGy916.16e3, MPMGy916.191c6, MPMGy916.195f6, MPMGy916.278b6, MPMGy916. |
| 751 | | | | d2rat1 | 0.0 | |
| 752 | 02 | 1 | RPCI32.406p4 | | | MPMGy916.113f2, MPMGy916.16e3, MPMGy916.191c6, MPMGy916.195f6, MPMGy916.278b6, MPMGy916. |
| 753 | 02 | 1 | RPCI31.49d1 | | | MPMGy916.113g2, MPMGy916.146d8, MPMGy916.172b7, MPMGy916.204a1, MPMGy916.271e12, MPMGy9 |
| 754 | 02 | 1 | RPCI32.423c7 | | | MPMGy916.112f1, MPMGy916.113f2, MPMGy916.146d8, MPMGy916.172b7, MPMGy916.186g6, MPMGy916 |
| 755 | 02 | 1 | RPCI31.168g3 | | | MPMGy916.112f1, MPMGy916.113g2, MPMGy916.146d8, MPMGy916.172b7, MPMGy916.204a1, MPMGy91 |
| 756 | 02 | 1 | RPCI32.368g3 | | | MPMGy916.113g2, MPMGy916.135d3, MPMGy916.172b7, MPMGy916.219g3, MPMGy916.265b11, MPMGy5 |
| 757 | 02 | 1 | RPCI32.420b8 | | | MPMGy916.113g2, MPMGy916.135d3, MPMGy916.172b7, MPMGy916.219g3, MPMGy916.265b11, MPMGy5 |
| 758 | 02 | 1 | RPCI32.420c8 | | | MPMGy916.113g2, MPMGy916.135d3, MPMGy916.172b7, MPMGy916.219g3, MPMGy916.265b11, MPMGy5 |

· 1 BAC associated with many YACs

· want to extract the list of YACs associated with each BAC

· BACa -> YAC1,YAC2,YAC3,…,YACm
· BACb -> YAC2,YAC3,YAC5,…,YACn

· examine linking relationships

A BAC

SOME YACs

RELATIONSHIP BETWEEN OUR DATA

# LWP::Simple

· It's very easy to grab a remote web page.

```perl
use LWP::Simple;
my $url  = "http://www.mdc-berlin.de/ratgenome/data/MDC-Map-15.html";
my $html = get($url);
```

· `$html` now contains the HTML content of the web page

```
HTML><HEAD><TITLE>MDC-Rat-Data</TITLE></HEAD>
<BODY scroll=yes>
<H1>Physical Mapping Data, Nov/01/2002</H1>
<P>Download: <A href="http://flipper.molgen.mpg.de:10085/mdcRATdata/MDC-
RatDataSet.tsv">MDC-RatDataSet.tsv</A>   (TAB separated values,
including RH-vectors, 1.8 MB)</P><BR>
<H3><U>Legend:</U></H3>
<TABLE border=0>
<TBODY>
<TR>
<TD><B>No </B></TD>
<TD>- consecutive number<BR></TD></TR>
<TR>
```

BIOINFORMATICS
Perl Workshop

TWO PROBLEMS – PART II

GENOME
SCIENCES
CENTRE

# Parsing HTML – HTML::TreeBuilder

- <u>Never</u> parse HTML with your own code, unless you have a good reason. Use existing parser modules.

```
use HTML::TreeBuilder;
my $tree = HTML::TreeBuilder->new_from_content($html);
```

- `$tree` is an object which you can traverse

- you have to know what you're looking for

# Examine HTML – Brittle!

```
<TABLE border=0>
<TBODY>
<TR>
<TD><B>No </B></TD>
<TD>- consecutive number<BR></TD></TR>
<TR>
<TD><B>Chr </B></TD>
<TD>- chromosome</TD></TR>
…
</TD></TR></TBODY></TABLE>
```

```
…
<TABLE rules=none border=1><FONT size=-1>
…
<TR bgColor=#eeeeee>
<TD> 748</TD>
<TD> 02</TD>
<TD> 1</TD>
<TD> RPCI31.64l18</TD>
<TD> </TD>
<TD> </TD>
<TD> MPMGy916.186d9, MPMGy916.34f11…
```



**Physical Mapping Data, Nov/01/2002**

# Fetch Columns from Second Table

Columns 2, 3, 6 contain data we want. Extract data and save in memory.

grep(?,@x)

IDIOM

```perl
# fetch table
my ($table) = grep($_->attr("rules") eq "none", $tree->find_by_tag_name("table"));
# get all rows from table
my @rows     = $table->find_by_tag_name("tr");
# for each row…
ROW:
foreach my $row (@rows) {
  # get all columns
  my @cols = $row->find_by_tag_name("td");
  # some columns do not contain data we want
  next unless @cols == 7;
  # get data from columns 2,3,6
  my $contig   = $cols[2]->as_text;
  my $bacname  = $cols[3]->as_textl
  my $yacnames = $cols[6]->as_text;
  # split YAC names a,b,c,d -> (a b c d)
  my @yacnames = split(/,/,$yacnames);
  # save data in a hash of lists
  push ( @{$bac_to_yacs{$bacname}}, @yacnames );
}
```
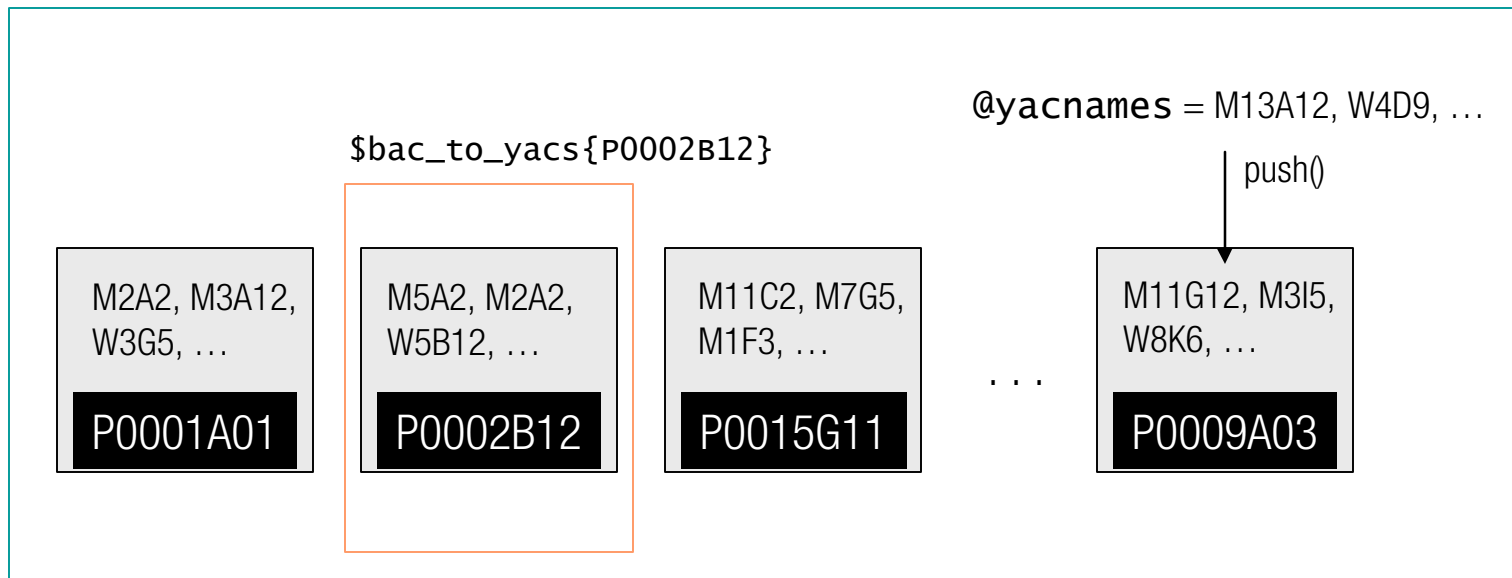
## Hashes and Arrays

```
my $bacname  = $cols[3]->as_text
my $yacnames  = $cols[6]->as_text;
my @yacnames = split(/,/,$yacnames);

push ( @{$bac_to_yacs{$bacname}}, @yacnames );
```

%bac_to_yacs



@yacnames = M13A12, W4D9, …

$bac_to_yacs{P0002B12}

push()

M2A2, M3A12, W3G5, …

P0001A01

M5A2, M2A2, W5B12, …

P0002B12

M11C2, M7G5, M1F3, …

P0015G11

. . .

M11G12, M3I5, W8K6, …

P0009A03

## Step 2 – Munge - Perl Makes It Easy

Store data in a way that allows you to easily find needed relationships – <u>choose wisely</u>

- · BAC -> list all associated YACs
  - · `@list = @{$bac_to_yac{$bacname}}`

- · BAC -> how many YACs?
  - · `scalar ( @list )`

- · how many total BACs?
  - · `scalar ( keys %bac_to_yac )`

- · how many total YACs?
  - · `$num_yacs = scalar ( map { @{$bac_to_yac{$_}} keys %bac_to_yac )`
  - · this sum doesn't take care of duplicates

- · how many average YACs per BAC?
  - · `use Math::VecStat qw(average);`
  - · `average ( map { scalar ( @{$bac_to_yac{$_}} ) } keys %bac_to_yac );`

# CPAN

- <u>CPAN</u> contains 5,000+ modules of all types – fun & serious
  - Perl Data Language (<u>PDL</u>) for matrix manipulation (PDL)
  - convert time to Swedish Chef speak (Acme::Time::Baby)

```
#!/usr/local/bin/perl
use Acme::Time::Baby language => "swedish chef";
print babytime "5:35";


Zee beeg hund is un zee sefen und zee little hund
is un zee six. Bork, bork, bork!
```

**SEARCH.CPAN.ORG**

- <u>Graph::Base</u> to create directed and undirected graphs

- <u>GraphViz</u> to generate GIF/TXT/EPS/PNG/…s from graph

# Standardized Module Documentation



STRING::RANDOM



MATH::VECSTAT

**NAME**

Grinder – grinds coffee

**SYNOPSIS**

use Grinder;

$g = Grinder->new();

$g->grind("coarse");

$g->empty();

**DESCRIPTION**

Models a Rancillio burr coffee grinder
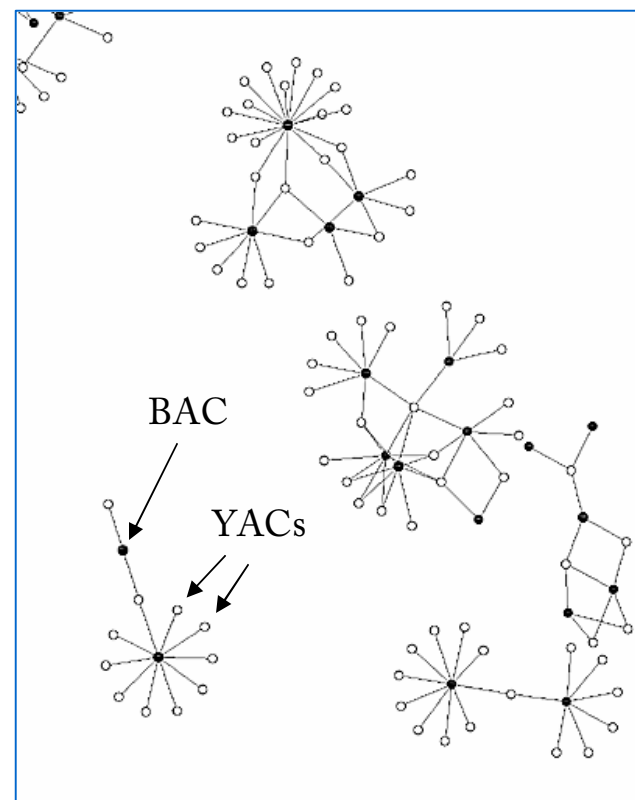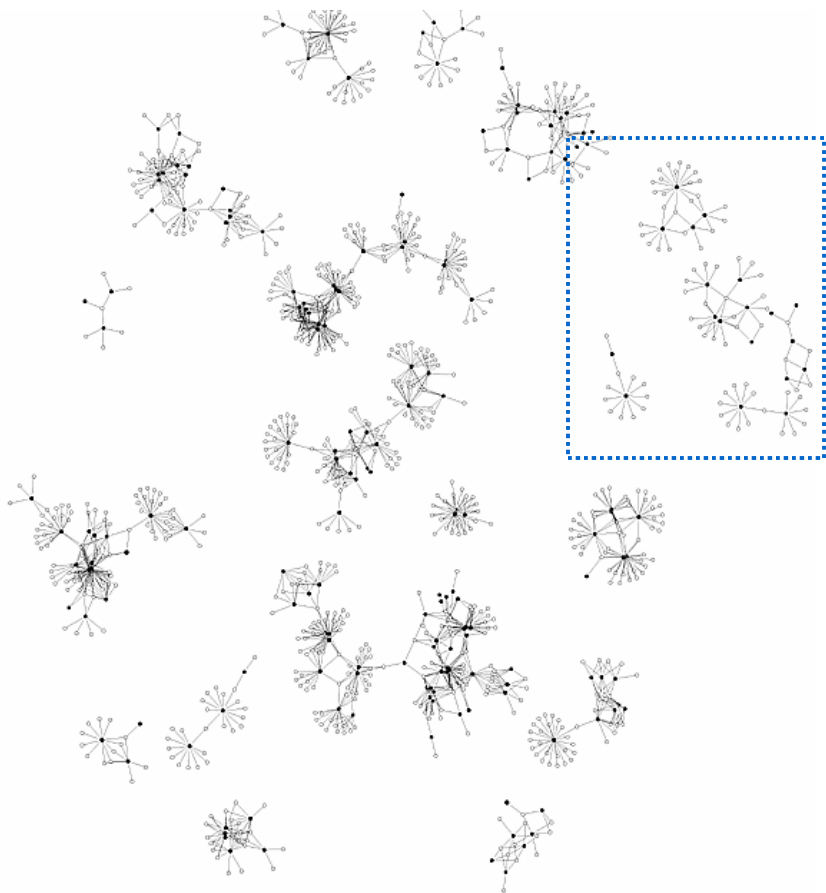
**HISTORY**

9 October 2003 - docs

**BUGS**

If found, remove from grinder

**AUTHOR**

M Krzywinski

# GraphViz – Big Bang for Little Buck



BAC

YACs

# Creating Graphs with Graph:: and GraphViz

map {} @x

IDIOM

BAD
USAGE

```perl
my $graph    = Graph::Undirected->new();
my $graphviz = GraphViz->new(directed=>0);

# for each BAC in the hash
foreach my $bac (keys %bac_to_yacs) {
    # get a list of all YACs for this BAC
    my @yacs = @{$bac_to_yacs{$bac}};
    # add edge between bac & yac in Graph::Undirected object
    map {$graph->add_edge($bac,$_) } @yacs;
    # for vizualization do the same for GraphViz object
    map { $graphviz->add_edge($bac,$_) } @yacs; # map {} IDIOM
    }
}

# create PNG image of graph
open(GRAPH,">/home/martink/www/htdocs/tmp/bacyac.png");
print GRAPH $graphviz->as_png;
close(GRAPH);
```
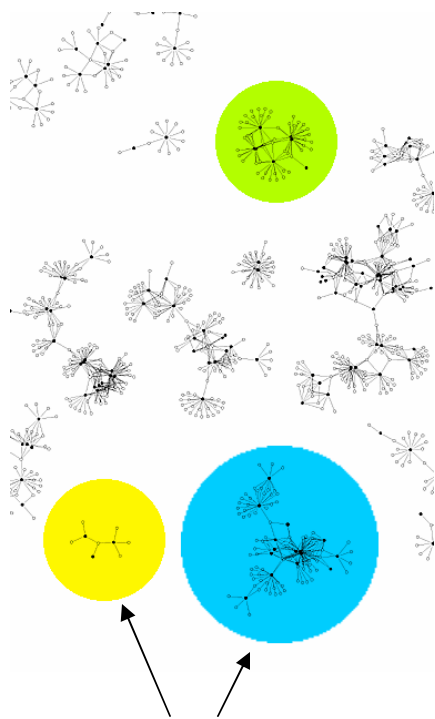
BIOINFORMATICS
Perl Workshop

TWO PROBLEMS – PART II

GENOME
SCIENCES
CENTRE

# List Clones in Contigs

List <u>connected components</u>, or contigs, created by BAC-YAC links.



```perl
# make a list of lists which contain connected vertices
my @groups = $graph->strongly_connected_components;
# iterate through each vertex list
foreach my $group_idx (0..@groups-1) {
    # get the vertices for this list
    my @vertices = @{$groups[$group_idx]};
    # for each vertex, report the group (contig) index,
    # vertex index and name
    foreach my $vertex_idx (0..@vertices-1) {
        printf("%d %d %s\n",
                $group_idx,
                $vertex_idx,
                $vertices[$vertex_idx]);
    }
}
```

contig is a
<u>connected component</u>

## Output - Create Output to STDOUT

It's nice to create output to STDOUT, rather than a file, because you can pipe your script into other processes.

```
foreach my $vertex_idx (0..@vertices-1) {
    printf("%d %d %s\n",
            $group_idx,
            $vertex_idx,
            $vertices[$vertex_idx]);
}
```

```
0 0 RPCI31.80h3
0 1 MPMGy916.380h8
0 2 WIBRy933.259d6
...
0 45 WIBRy933.284b4
0 46 WIBRy933.219c12
0 47 MPMGy916.110d1
1 0 MPMGy916.369g12
1 1 MPMGy916.282g2
1 2 RPCI31.33m10
...
21 7 RPCI31.17n14
21 8 WIBRy933.106h8
21 9 RPCI31.17i17
```

· Perl is friendly – you can copy file handles
  · STDOUT to file
  · file to STDOUT

CONTIG    CONTIG    CLONE NAME

CLONE INDEX

# Munge at Prompt

Don't forget that the command prompt offers powerful tools to manipulate and extract data – <u>generate maximally detailed reports</u> and parse later

```
0 0 RPCI31.80h3
0 1 MPMGy916.380h8
0 2 WIBRy933.259d6
...
0 45 WIBRy933.284b4
0 46 WIBRy933.219c12
0 47 MPMGy916.110d1
1 0 MPMGy916.369g12
1 1 MPMGy916.282g2
1 2 RPCI31.33m10
...
21 7 RPCI31.17n14
21 8 WIBRy933.106h8
21 9 RPCI31.17i17
```

· how many contigs?
  · cut –d " " –f 1 data.txt | sort –u | wc

· how many clones?
  · cut –d " " –f 3 data.txt | sort –u | wc

· how many clones in contig 10?
  · grep –d "^10 " data.txt | wc

· which contigs have < 20 clones?
  · cut –d " " –f 1 data.txt | uniq –c | egrep " 1?[0-9] "

| clones | contig | clones | contig |
|--------|--------|--------|--------|
| 16 | 13 | 11 | 18 |
| 18 | 14 | 8 | 19 |
| 18 | 15 | 9 | 20 |
| 13 | 16 | 10 | 21 |

# Perl

PRODUCTIVE

CREATIVE

LINGUAL

COMPACT

OPEN SOURCE

DOES NOT SPIT