

BioRuby + KEGG API + KEGG DAS = wiring knowledge for genome and pathway

Toshiaki Katayama <k@bioruby.org>
Human Genome Center, University of Tokyo, Japan
<http://bioruby.org/>
<http://www.genome.jp/kegg/soap/>
<http://das.hgc.jp/>

What is BioRuby?

- Yet another BioPerl written in Ruby
 - since Nov 2000
- Mainly developed in Japan
 - including supports for Japanese resources like KEGG

What is Ruby?

- Created by Japanese author 'matz'
- Scripting language
 - clean syntax, easy to learn, powerful enough
- Purely object oriented
 - Integer, String, Regexp, Exception etc. w/o exception
- Sufficient libraries
 - You can use most of the BioRuby functionality without install additional libraries.



© 2004 Toshiaki Katayama, Human Genome Center, University of Tokyo

BioRuby.org

What you can do with BioRuby?

- Biological sequence manipulation
- Database entry retrieval and parsing
- Running and parsing usual applications
- Graph computation

© 2004 Toshiaki Katayama, Human Genome Center, University of Tokyo

BioRuby.org

Sequence manipulations

- Splicing
- Reverse complement
- Translation
- Composition
- Molecular weight
- Window search

```
# print in FASTA format
puts seq.to_fasta("foo", 60)
```

```
# for selenoproteins
ct =Bio::CodonTable.copy(1)
ct['tga'] = 'U'

puts seq.translate
puts seq.translate(1, ct)
```

```
#!/usr/bin/env ruby

require 'bio'

seq = Bio::Sequence::NA.new(ARGF.read)

puts seq.subseq(1,3)
puts seq.splicing("join(1..23,45..67)")
puts seq.complement
puts seq.translate
puts seq.gc_percent
puts seq.composition

seq.window_search(15, 3) do |subseq|
  peptide = subseq.translate
  puts peptide.molecular_weight
end
```

© 2004 Toshiaki Katayama, Human Genome Center, University of Tokyo

BioRuby.org

Entry parsers

- GenBank, GenPept, RefSeq, DDBJ
- EMBL, UniProt (TrEMBL, SwissProt)
- PDB
- KEGG/GenomeNet
 - GENES, GENOME, ENZYME, COMPOUND, KO, BRITE, CELL, Expression, Keggtab, AAindex
- GFF
- GO
- FANTOM
- Transfac, Prosite
- LITDB, MEDLINE
- NBRF, PIR
- FASTA format

```
#!/usr/bin/env ruby
# Usage:
# % auto.rb dbfile

require 'bio'

Bio::FlatFile.auto(ARGF) do |ff|
  ff.each do |entry|
    # do something
  end
end
```

© 2004 Toshiaki Katayama, Human Genome Center, University of Tokyo

BioRuby.org

Database access

- OBDA (Open Bio* Database Access)
 - BioRegistry
 - BioFlat
 - BioFetch
 - BioSQL
- PubMed
- DAS
 - Ensembl, WormBase etc.
 - KEGG DAS
- SOAP
 - KEGG API
 - DDBJ XML
 - NCBI ESOAP

© 2004 Toshiaki Katayama, Human Genome Center, University of Tokyo



OBDA configuration

- ~/.bioinformatics/seqdatabase.ini
- /etc/bioinformatics/seqdatabase.ini
- <http://www.open-bio.org/registry/seqdatabase.ini>

```
VERSION=1.00

[genbank]
protocol=flat
location=/export/database/
dbname=genbank

[swissprot]
protocol=biosql
location=db.bioruby.org
dbname=biosql
driver=mysql
biodbname=sp

[embl]
protocol=biofetch
location=http://bioruby.org/cgi-bin/biofetch.rb
dbname=embl
```

```
#!/usr/bin/env ruby
require 'bio'
reg = Bio::Registry.new
sp = reg.get_database('swissprot')
puts sp.get_by_id('CYC_BOVIN')
gb = reg.get_database('genbank')
puts gb.get_by_id('AA2CG')
```

© 2004 Toshiaki Katayama, Human Genome Center, University of Tokyo



Running applications

- Bio::Blast, Fasta, HMMER, EMBOSS
- Bio::ClustalW, MAFFT
- Bio::Genscan
- Bio::PSORT, TargetP
- Bio::SOSUI, TMHMM

```
#!/usr/bin/env ruby
require 'bio'

File.open("my_blast_output.xml") do |file|
  Bio::Blast.reports(file) do |report|
    report.hits do |hit|
      hit.each do |hsp|
        puts hsp.query_id, hsp.target_id,
              hsp.bit_score, hsp.evaluate, hsp.overlap
      end
    end
  end
end
```

© 2004 Toshiaki Katayama, Human Genome Center, University of Tokyo

Bioinformatics
BioRuby.org

What is KEGG? Kyoto Encyclopedia of Genes and Genomes

- <http://www.genome.jp/kegg/>
 - PATHWAY
 - LIGAND
 - GENES, GENOME
 - KO
 - SSDB
 - Expression
 - Glycan
 - etc.

KEGG: Kyoto Encyclopedia of Genes and Genomes

A grand challenge in the post-genomic era is a complete computer representation of the cell and the organism, which will enable computational prediction of higher-level complexity of cellular processes and organism behaviors from genomic information. Towards this end we have been developing a bioinformatics resource named KEGG, Kyoto Encyclopedia of Genes and Genomes, as part of the research projects in the Kanehisa Laboratory of Kyoto University Bioinformatics Center.

Building blocks of life		Wiring diagrams of life	
Genetic materials	Chemical materials	Genetic networks	Chemical networks
KEGG Gene Universe Genes and gene products in complete genomes	KEGG Chemical Universe Chemical compounds Complex carbohydrates Active peptides	KEGG Protein Network Metabolic pathways Regulatory pathways Molecular complexes	
Genomic contexts Orthology/paralog relations KO (KEGG Orthology)	Chemical reactions RC (Reaction Classification)	Network-network relations Network-environment relations Diseases	

Genetic and chemical blueprints of life
 KEGG Table of Contents
 KEGG Release 30.0, April 2004 (plus daily updates)

© 2004 Toshiaki Katayama

KEGG PATHWAY

KEGG PATHWAY Database

Current knowledge on molecular interaction networks, including metabolic pathways, regulatory pathways, and molecular complexes

Go to: **1. Metabolism**
2. Genetic Information Processing
3. Environmental Information Processing
4. Cellular Processes
5. Human Diseases

See also: **KO (KEGG Orthology)**

1. Metabolism

1.1 Carbohydrate Metabolism

- Glycolysis / Gluconeogenesis Ortholog, Oxidoreductases
- Citrate cycle (TCA cycles) Ortholog
- Pentose phosphate pathway Ortholog
- Pentose and glucuronate interconversions Ortholog
- Fructose and mannose metabolism Ortholog, PTS
- Galactose metabolism Ortholog
- Ascorbate and aldarate metabolism Ortholog
- Pyruvate metabolism Ortholog
- Glyoxylate and dicarboxylate metabolism Ortholog
- Propanoate metabolism Ortholog
- Butanoate metabolism Ortholog
- CS-Branch eddic acid metabolism Ortholog
- Inositol metabolism Ortholog

1.2 Energy Metabolism

- Oxidative phosphorylation Ortholog
- ATP synthesis Ortholog
- Photosynthesis Ortholog
- Carbon fixation Ortholog
- Reductive carboxylate cycle (CO2 fixation) Ortholog
- Methane metabolism Ortholog
- Nitrogen metabolism Ortholog
- Sulfur metabolism Ortholog

1.3 Lipid Metabolism

- Fatty acid biosynthesis (path 1) Ortholog
- Fatty acid biosynthesis (path 2) Ortholog

Glycolysis / Gluconeogenesis - Escherichia coli K-12 MG1655

Go to: [LinkDB search](#) | [Orthology Table](#)

Go to: [Escherichia coli K-12 MG1655](#) | [Exec](#)

KEGG GENES

DBGET Result: E.coli b0002

[LinkDB](#) | [SSDB](#) | [FASTA-genes](#) | [FASTA-sp](#) | [BLAST-nr](#) | [SOSUI](#) | [SPORT](#) | [MOTIF](#)

ENTRY: **b0002** CDS: **E.coli**

NAME: thrA, thrA1, thrA2

DEFINITION: aspartokinase I / homoserine dehydrogenase I [EC:2.7.2.4.1.1.1.3]

ORTHODOLOGY: RO: [K00003](#) homoserine dehydrogenase

KO: [K00028](#) aspartate Kinase

CLASS: Metabolism; Amino Acid Metabolism; Glycine, serine and threonine metabolism [PATH:[eco00260](#)]

Metabolism; Amino Acid Metabolism; Lysine biosynthesis [PATH:[eco00300](#)]

POSITION: 337..2799

DBLINKS: Wisconsin: [b0002](#)
Colibri: [thrA](#)
NCBI: [1786183](#)
SPT: [P00561](#)

CODON_USAGE

	T	C	A	G
T	11	19	10	13
C	8	13	2	43
A	30	15	1	23
G	19	18	5	27

AASeq: 820

NRVLEFGTGSVANAERFLRVADILESNRQGGVAVLSAPAKITNHLVAMIKTTISGDDA
LFINISDAERFIALLTGLAAQGFPLAQLKTFVDFQFAIKHVLHGILSLGGCPDSINA
ALICRGERKMSIALMAGVLEAGHNVIVDFVEKLVAGHYLESTVDIAESTRRIAASRIP
ADHMLMAGFTAGNEKGEVLVLRNGSDVSAVLAACLRAADCEIIVTDVGVYTCDFPROV
PDRARLKSMSYQEMELSYGAKVLFHRTITPFAQFIPCLIKNTGNPQAPGTLIGASRD
EDELFRKISLNLANMNFVSGFGRKHVQAAARVFAVMSRARSVLLTQSSSERSISF
CVFQCVRAERMQEFTVLEKREGLLEPLAVTERLAHISVVGDMRTLRGISARFPAAL
ARANINIVATAGSSERSISVVNNDATTGVRVTHQMLFNTDQVIEVFVIGVGGVGGAL
LEQLRQSQVLEKNIHDLKRVCGVANSKALLTVHGLNLENWQELAQAFEPNLRGLRLR
VKYEHLLNPIVVDCTSQAQVADQYADFLREGFHVTVPNKANTSSMDYHQLRYAAEKSR
RRFLYDTHVAGLPIENLQNLNAGDELKMSLISGLSLYIFKGLDEGNSFSEATTLA
REMGVTEPRDLDGSDVARKLLILARETGRELEADTEFVFAEFAEFAEFAEFAEFAEFAE
NLSQDDLFARVAKAREGKRVLYRQNIDEDGCRVKIAEVDKNDPLFRKNGENALAF
YSHYQPLPLVLRGYAGNDVTAAGVFADLLRLTSLKGLV
2463

atgcgagtgtagaattcggcggatcacatcggcgaatgcagaacgtttctcgctgtt
gcgatattctggaagcaatgccagcaggggcagtgccaccgctctctcgccoc
gccaaatcaccaccaccctggtgagatgtaaaaccatagcggccagatgct
taaccatcacggatgccagcagattttggcgaatttgacggagatcgccgc
gccagcgggttcccgctggcgaattgaaaacttctgcatcaggaattgccaa
ataaacatgctcgtgcatgattgatttggggcagtgcccgatgcatcaacgct
gcgctgatttgcgtagcagaaatgctgcatgcaattgcccggcgtattagaagc
cggctcaccacgtactgttatcgatccggcgaataactgctggcagtgggcattac

LinkDB Search Result

Database: LinkDB

Database of Link Information
Release 04-04-16, Apr 04
Institute for Chemical Research, Kyoto University
110,095,443 entries

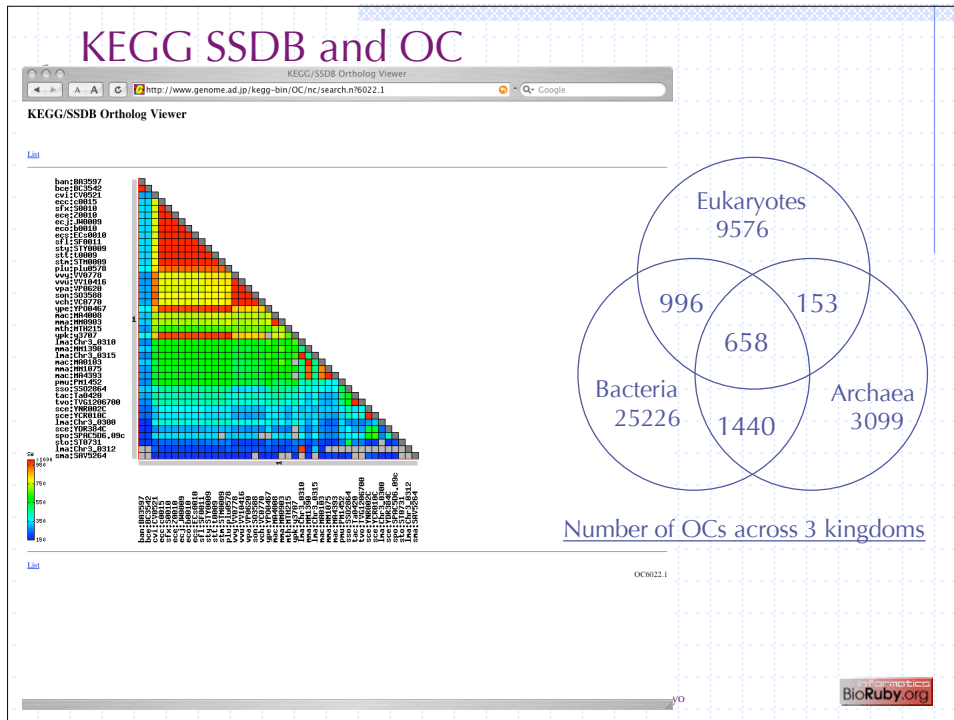
From: E.coli:b0002
To: All

291 hits from 12+ databases

1. [GenBank \(159\)](#)
2. [EMBL \(22\)](#)
3. [PIR \(50\)](#)
4. [PRF \(25\)](#)
5. [SPT/EMBL \(1\)](#)
6. [PDB \(3\)](#)
7. [ECHOSEA \(8\)](#)
8. [KO \(2\)](#)
9. [COMPOUND \(1\)](#)
10. [ENZYME \(2\)](#)
11. [PATHWAY \(4\)](#)
12. [GENOME \(1\)](#)
13. [OTHERS \(3\)](#)
14. [All databases \(291\)](#)

[Link table for E.coli](#)

DBGET integrated database retrieval system, [GenomeNet](#)



- ## KEGG standardization
- KGML
 - ISMB2004 poster C-45
 - KEGG API
 - ISMB2004 poster G-17
 - KEGG DAS
 - ISMB2004 poster E-29
- © 2004 Toshiaki Katayama, Human Genome Center, University of Tokyo
- BioRuby.org

KGML

KGML (KEGG Markup Language)

The KEGG Markup Language (KGML) is an exchange format of the KEGG graph objects maps that are manually drawn and updated. KGML enables automatic drawing of KEGG facilities for computational analysis and modeling of protein networks and chemical net metabolic pathways contain two types of graph objects, how boxes (enzymes) are link (compounds) are linked by a reaction in the KEGG pathway diagrams. In contrast, the pathways contain only the aspect of how boxes (proteins) are linked by a relation. Note (KEGG Orthology) identifies in the current KEGG system, but for historical reasons so are marked with EC numbers in the actual pathway diagrams.

Documents

- KEGG Markup Language manual
- KGML v0.4 DTD [dtd | html]
- KGML v0.4 Readme [txt | html]

Data


- KEGG reference metabolic pathways (Last update: Apr 23, 2004)
- KEGG reference regulatory pathways (Last update: Apr 21, 2004)
- KEGG metabolic pathways linked to KO (Last update: Apr 23, 2004)
- KEGG regulatory pathways linked to KO (Last update: Mar 25, 2004)
- KEGG organism-specific metabolic pathways
 - Select organism:
- KEGG organism-specific regulatory pathways
 - Select organism:

Examples

- Pathway viewer using KGML [map00770 | eco00770 | hsa00770]


Previous Versions

```
<?xml version="1.0"?>
<!DOCTYPE pathway SYSTEM "http://www.genome.ad.jp/kegg/xml/KGML_v0.4_dtd">
<!-- Creation date: Apr 01 2004 01:22:17 +0900 (JST) -->
<pathway name="pathmap00010" org="map" number="00010"
  title="Glycolysis / Gluconeogenesis"
  image="http://www.genome.ad.jp/kegg/pathway/map/map00010.gif"
  link="http://www.genome.ad.jp/dbget-bin/show_pathway?map00010">
  <entry id="1" name="ec1.2.1.3" type="enzyme" reaction="rnrR00710">
    link="http://www.genome.ad.jp/dbget-bin/www_bget?enzyme=1.2.1.3">
    <graphics name="1.2.1.3" fgcolor="#000000" bgcolor="#FFFFFF"
      type="rectangle" x="170" y="1018" width="45" height="17"/>
  </entry>
  <entry id="2" name="ec6.2.1.1" type="enzyme" reaction="rnrR00235">
    link="http://www.genome.ad.jp/dbget-bin/www_bget?enzyme=6.2.1.1">
    <graphics name="6.2.1.1" fgcolor="#000000" bgcolor="#FFFFFF"
      type="rectangle" x="102" y="916" width="46" height="17"/>
  </entry>
  <entry id="3" name="ec1.2.1.5" type="enzyme" reaction="rnrR00711">
    link="http://www.genome.ad.jp/dbget-bin/www_bget?enzyme=1.2.1.5">
    <graphics name="1.2.1.5" fgcolor="#000000" bgcolor="#FFFFFF"
      type="rectangle" x="170" y="1039" width="45" height="17"/>
  </entry>
  <entry id="4" name="cpd:C00033" type="compound">
    link="http://www.genome.ad.jp/dbget-bin/www_bget?compound=C00033">
    <graphics name="C00033" fgcolor="#000000" bgcolor="#FFFFFF"
      type="circle" x="102" y="971" width="8" height="8"/>
  </entry>
  <entry id="5" name="pathmap00650" type="map">
    link="http://www.genome.ad.jp/kegg/pathway/map/map00650.html">
    <graphics name="Butanoate metabolism" fgcolor="#000000" bgcolor="#FFFFFF"
      type="roundrectangle" x="645" y="924" width="128" height="25"/>
  </entry>
  <entry id="6" name="pathmap00660" type="map">
    link="http://www.genome.ad.jp/kegg/pathway/map/map00660.html">
    <graphics name="C5-Branched dibasic acid metabolism" fgcolor="#000000" bgcolor="#FFFFFF"
      type="roundrectangle" x="637" y="898" width="205" height="25"/>
  </entry>
  <entry id="7" name="pathmap00640" type="map">
    link="http://www.genome.ad.jp/kegg/pathway/map/map00640.html">
    <graphics name="Propanoate metabolism" fgcolor="#000000" bgcolor="#FFFFFF"
      type="roundrectangle" x="637" y="864" width="131" height="25"/>
  </entry>
  <entry id="8" name="pathmap00710" type="map">
    link="http://www.genome.ad.jp/kegg/pathway/map/map00710.html">
    <graphics name="Carbon fixation" fgcolor="#000000" bgcolor="#FFFFFF"
      type="roundrectangle" x="643" y="720" width="90" height="25"/>
  </entry>
```

© 2004 Toshiaki Katayama, Human Genome Center, University of Tokyo 

KEGG API

- <http://www.genome.jp/kegg/soap/>
- SOAP/WSDL based web service
 - XML, HTTP
 - Can be accessed by any language
- Proteome and pathway analysis
 - KEGG GENES, SSDB, PATHWAY
 - DBGET, LinkDB
- Updated to v3.0 (May 2004)

© 2004 Toshiaki Katayama, Human Genome Center, University of Tokyo 

Example: obtain homologs and motifs by SSDB

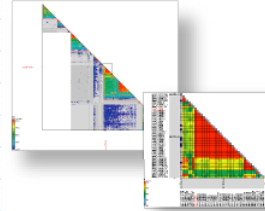
```
#!/usr/bin/env ruby

require 'bio'

serv = Bio::KEGG::API.new

homologs = serv.get_all_best_neighbors_by_gene("hsa:7368")

homologs.each do |hit|
  gene = hit.genes_id2
  if motifs = serv.get_motifs_by_gene(gene, "pfam")
    motifs.each do |motif|
      name = motif.motif_id
      desc = motif.definition
      puts "#{gene}: #{name} #{desc}"
    end
  end
end
```



© 2004 Toshiaki Katayama, Human Genome Center, University of Tokyo

BioRuby.org

Example: KO, OC, PC

```
#!/usr/bin/env ruby

require 'bio'

serv = Bio::KEGG::API.new

ko_list = serv.get_ko_by_gene("eco:b0002")

# list of genes assigned same KO (KEGG orthology)
ko_genes = serv.get_ko_members(ko_list.first)

# list of genes assigned to same OC (ortholog cluster)
oc_genes = serv.get_all_oc_members_by_gene("hsa:7368")

# list of genes assigned to same PC (paralog cluster)
pc_genes = serv.get_all_pc_members_by_gene("hsa:7368")

puts "# KO", ko_genes, "# OC", oc_genes, "# PC", pc_genes
```

© 2004 Toshiaki Katayama, Human Genome Center, University of Tokyo

BioRuby.org

Example: getting database info/entries

```
#!/usr/bin/env ruby

require 'bio'

serv = Bio::KEGG::API.new

# lists available organisms in KEGG
orgs = serv.list_organisms
orgs.each do |entry|
  puts "#{entry.entry_id} #{entry.definition}"
end

# lists available pathways in KEGG
list = serv.list_pathways("hsa")
list.each do |entry|
  puts "#{entry.entry_id} #{entry.definition}"
end

# getting EMBL entry
puts serv.bget("embl:BUM")

# getting entries from KEGG GENES
puts serv.bget("hsa:7368 hsa:7369")
```

© 2004 Toshiaki Katayama, Human Genome Center, University of Tokyo

BioRuby.org

Example: using PATHWAY

```
#!/usr/bin/env ruby

require 'bio'

serv = Bio::KEGG::API.new

# lists gene ids on pathway
genes = serv.get_genes_by_pathway("path:hsa00020")
puts "# genes on human's pathway 00020"
genes.each do |gene|
  puts gene
end

# converts EC numbers to genes
list = ["ec:1.1.1.1", "ec:1.2.1.1"]
list.each do |ec|
  puts "# E. coli genes for #{ec}"
  puts serv.get_genes_by_enzyme(ec, "eco")
end
```

© 2004 Toshiaki Katayama, Human Genome Center, University of Tokyo

BioRuby.org

Example: coloring PATHWAY

```
#!/usr/bin/env ruby

require 'bio'

serv = Bio::KEGG::API.new

# mark pathway
objs = ['eco:b0002', 'cpd:C00263']
url1 = serv.mark_pathway_by_objects('path:eco00260', objs)

# color pathway
fg_list = ['blue', 'green']
bg_list = ['#ff0000', 'yellow']
url2 = serv.color_pathway_by_objects('path:eco00260',
    objs, fg_list, bg_list)

# save the result images
serv.save_image(url1, "marked_pathway.gif")
serv.save_image(url2, "colored_pathway.gif")
```

© 2004 Toshiaki Katayama, Human Genome Center, University of Tokyo



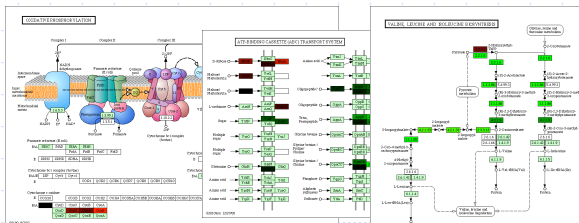
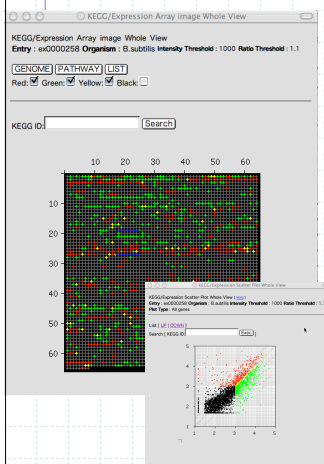
Example: gene expression and pathway analysis

```
serv = Bio::KEGG::API.new

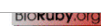
list = serv.get_genes_by_pathway("path:bsu00020")
fg_colors = Array.new
bg_colors = Array.new

list.each do |gene|
  fg_colors << "black"
  bg_colors << ratio2rgb(gene)
end

url = serv.color_pathway_by_objects(
  "path:bsu00020", list, fg_colors, bg_colors)
```



© 2004 Toshiaki Katayama, Human Genome Center, University of Tokyo



Example: mapping PDB on pathway

```
#!/usr/bin/env ruby
require 'bio'

serv = Bio::KEGG::API.new

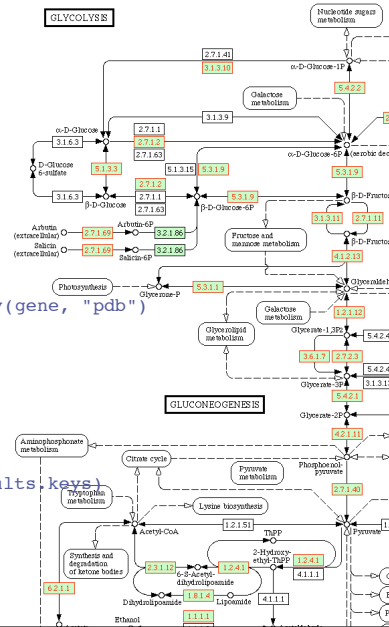
# list of genes on specified pathway
path = ARGV.shift || "path:eco00010"
genes = serv.get_genes_by_pathway(path)

# search DB links to PDB using LinkDB
results = Hash.new
genes.each do |gene|
  if pdb_links = serv.get_all_linkdb_by_entry(gene, "pdb")
    pdb_links.each do |link|
      results[gene] = true
    end
  end
end

# generates colored image
url = serv.mark_pathway_by_objects(path, results.keys)

# save the image
serv.save_image(url, "linked_to_pdb.gif")
```

© 2004 Toshiaki Katayama, Human Genome Centre



KEGG DAS

- <http://das.hgc.jp/>
 - Currently including 188 organisms
 - Data from KEGG GENOME and GENES
 - Build on to of the GMOD/GBrowse

© 2004 Toshiaki Katayama, Human Genome Center, University of Tokyo

BioRuby.org

KEGG/DAS server

http://das.hgc.jp/

KEGG DAS

GMOD and BioDAS server for KEGG GENOME, GENES and SSDB databases (test release).

- DAS dsn

```
[GBrowse] [DAS] KEGG aae : Aquifex aeolicus VF5
[GBrowse] [DAS] KEGG afu : Archaeoglobus fulgidus VC-16
[GBrowse] [DAS] KEGG ana : Anabaena sp. PCC 7120
[GBrowse] [DAS] KEGG ape : Aeropyrum pernix K1
[GBrowse] [DAS] KEGG atc : Agrobacterium tumefaciens C58
[GBrowse] [DAS] KEGG ath : Arabidopsis thaliana
[GBrowse] [DAS] KEGG atu : Agrobacterium tumefaciens C58
[GBrowse] [DAS] KEGG bab : Buchnera aphidicola Bp
[GBrowse] [DAS] KEGG ban : Bacillus anthracis Ames
[GBrowse] [DAS] KEGG bas : Buchnera aphidicola Sg
[GBrowse] [DAS] KEGG bba : Bdellovibrio bacteriovorus HD100
[GBrowse] [DAS] KEGG bbr : Bordetella bronchiseptica RB50
[GBrowse] [DAS] KEGG bbu : Borrelia burgdorferi B31
[GBrowse] [DAS] KEGG bca : Bacillus cereus ATCC 10987
[GBrowse] [DAS] KEGG bce : Bacillus cereus ATCC 14579
[GBrowse] [DAS] KEGG bfl : Blochmannia floridanus
[GBrowse] [DAS] KEGG bha : Bacillus halodurans C-125
[GBrowse] [DAS] KEGG bja : Bradyrhizobium japonicum USDA110
[GBrowse] [DAS] KEGG blo : Bifidobacterium longum NCC2705
[GBrowse] [DAS] KEGG bme : Brucella melitensis 16M
[GBrowse] [DAS] KEGG bms : Brucella suis 1330
[GBrowse] [DAS] KEGG bpa : Bordetella parapertussis 12822
[GBrowse] [DAS] KEGG bpe : Bordetella pertussis Tohama I
[GBrowse] [DAS] KEGG bsu : Bacillus subtilis 168
[GBrowse] [DAS] KEGG bth : Bacteroides thetaotaomicron VPI-5482
[GBrowse] [DAS] KEGG buc : Buchnera aphidicola APS
[GBrowse] [DAS] KEGG cac : Clostridium acetobutylicum ATCC 824
[GBrowse] [DAS] KEGG cbu : Coxiella burnetii RSA 493
[GBrowse] [DAS] KEGG cca : Chlamydia caviae GPIC
[GBrowse] [DAS] KEGG ccr : Caulobacter crescentus CB15
[GBrowse] [DAS] KEGG cdi : Corynebacterium diphtheriae gravis NCTC13129
[GBrowse] [DAS] KEGG cef : Corynebacterium efficiens YS-314
[GBrowse] [DAS] KEGG cel : Caenorhabditis elegans
[GBrowse] [DAS] KEGG cgl : Corynebacterium glutamicum ATCC 13032
[GBrowse] [DAS] KEGG cje : Campylobacter jejuni NCTC11168
[GBrowse] [DAS] KEGG cmu : Chlamydia muridarum
[GBrowse] [DAS] KEGG cna : Chlamydia pneumoniae AR39
```

KEGG ana : Anabaena sp. PCC 7120

Showing 20 kbp from ana, positions 2,364,845 to 2,384,844

Instructions: Search using a sequence name, gene name, locus, or other landmark. The wildcard character * is allowed. To center on a location, click the ruler. Use the Scroll/Zoom buttons to change magnification and position.

Examples: ana, alpha, beta, gamma, delta, epsilon, zeta.

[Hide banner] [Hide instructions] [Bookmark this view] [Link to an image of this view] [Publication quality image] [Help]

Landmark or Region: ana2364845..2384844

Search: Flip

Scroll/Zoom: Show 20 kbp

Overview of ana

CDS

Gene

air1978	air1979	air1980	air1981	air1982	air1983	air1984	rml65a	trnW-11e-1	air1986	rnc55a	air1988	air1987	air1989
---------	---------	---------	---------	---------	---------	---------	--------	------------	---------	--------	---------	---------	---------

Motif

ABC_membrane, ABC_tran, ABC_TRANSPORTER, ABC_TRANSPORTER, ABC_TRANSPORTER, Glycos_transf_2, BUF942, Glucosylase, BOK, Competence, gljD, trnL_cant_2e, rpl_TRNA_LIGASE_1L_2, Transposase_20, Transposase_9, TPRT, recI, Dik

3-frame translation (forward)

DMU/CG content

3-frame translation (reverse)

Pathway

ana00010	ana00260	ana00310	ana00310	ana00310
ana00092	ana00270			
ana00500				
ana00522				

EC

2.7.1.2	6.1.1.14	3.1.1.1
---------	----------	---------

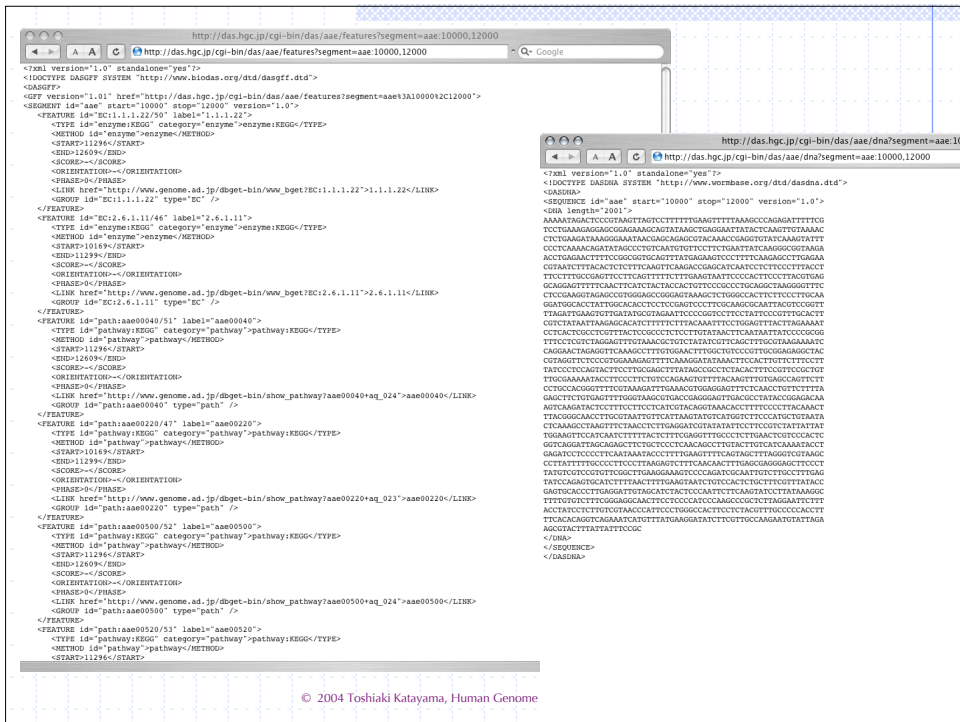
refSeq

ANE2944	ANE2128	ANE2935	ANE29031(ANE65020)	ANE29343	ANE26514	ANE1957	ANE29327	ANE19374
1,798,989,839	1,497,301,626	9,560,011,255	9,652,907,973	9,660,989,727	1,177,439,239	1,407,125,404	0,622,795,644	1,066,389,771

Data Source

KEGG ana - Anabaena sp. PCC 7120

Dumps, Searches and other Operations:



Example: sequence and annotations by DAS

```
#!/usr/bin/env ruby

require 'bio'

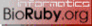
serv = Bio::DAS.new("http://das.hgc.jp/cgi-bin/")

segment = Bio::DAS::SEGMENT.region("I", 1001, 2000)

# get DNA sequence from S. cerevisiae
list = serv.get_dna("sce", segment)
list.each do |dna|
  puts dna.sequence
end

# get features
list = serv.get_features("sce", segment)
list.segments.each do |segment|
  segment.features.each do |feature|
    puts feature.entry_id
    puts feature.start
  end
end
```

© 2004 Toshiaki Katayama, Human Genome Center, University of Tokyo



Acknowledgements

- BioRuby developers
 - Naohisa Goto, Mitsuteru Nakao, Yoshinori Okuji, Shuichi Kawashima, Masumi Itoh, Alex Gutteridge and some other contributors on the net.
- KEGG curators and KEGG API developers
 - Bioinformatics center, Human genome center
 - Yoko Sato, Miho Matsubayashi, Satoshi Miyazaki
- KEGG DAS
 - Mayumi Takashio, Mari Watanabe
- Open Bio* community

© 2004 Toshiaki Katayama, Human Genome Center, University of Tokyo



URL

- BioRuby
 - <http://bioruby.org/>
 - CVS, ML - hosted at open-bio.org
- KEGG API
 - <http://www.genome.jp/kegg/soap/>
- KEGG DAS
 - <http://das.hgc.jp/>
- ISMB posters at E-29 G-17 C-45

© 2004 Toshiaki Katayama, Human Genome Center, University of Tokyo

