

JSTゲノムリテラシー講座
ゲノムアノテーション-データベースとツールの活用-

KEGG DASとKEGG APIで繋がる ゲノムアノテーションとパスウェイ

東京大学医科学研究所ヒトゲノム解析センター
片山俊明 <k@bioruby.org>
<http://bioruby.org/>
<http://www.genome.jp/kegg/soap/>
<http://das.hgc.jp/>

2004/08/24@サイエンスプラザ

目次

- KEGG (Kyoto Encyclopedia of Genes and Genomes)
 - GENES
 - KO (KEGG Orthology)
 - SSDB
 - OC (Ortholog Cluster)
 - PATHWAY
- KEGG DAS
 - DAS (Distributed Annotation System)
- KEGG API
 - SOAP/WSDL
- BioRuby
 - Ruby (Object oriented scripting language)

KEGG

- <http://www.genome.jp/kegg/>
- 京都大学化学研究所バイオインフォマティクスセンター(金久研)で開発されているデータベース
 - PATHWAY
 - LIGAND
 - GENES, GENOME
 - KO
 - SSDB
 - Expression
 - Glycan
 - etc.

© 2004 Toshiaki Katayama,



KEGGのゲノムアノテーション

- GENES - protein universe
- ゲノムの決まった生物種(現在約200種)の全遺伝子

DBGET Result: E.coli b0002

Escherichia coli K-12 MG1655: b0002

Entry: b0002 CDS E.coli

Gene name: thrA, thrA1, thrA2

Definition: aspartokinase I / homoserine dehydrogenase I [EC:2.7.2.4 1.1.1.3]

KO: K01003 homoserine dehydrogenase, K01095 aspartate kinase

Class: Metabolism; Amino Acid Metabolism; Glycine, serine and threonine metabolism [PATH:ec00260], Metabolism; Amino Acid Metabolism; Lysine biosynthesis [PATH:ec00300]

SSDB: Ortholog, Paralog, Motif, Gene cluster

Other DBs: Wisconsin: b0002, Colibri: thrA, RegulonDB: ECK12000002, NCBI-cgi: 1781813, NCBI-geneID: 945803, UniProt: P09561

Position: 337..2799 Genome map

AA seq: 820 aa (AA seq, FASTA-genes, FASTA-sp, BLAST-nr)

NT seq: 2463 nt (NT seq, +upstream 0 nt, +downstream 0 nt)

Motifs in eco:b0002: Aspartate kinase, Homoserine dehydrogenase

KEGG 550B Gene Cluster Search: 550B ID: eco0002 (2014-01-01) (2017), Definition: aspartokinase I / homoserine dehydrogenase I, Gene: thrA, thrA1, thrA2, Gene size: 2463 bp

KEGG 550B Gene Cluster Search: 550B ID: eco0002 (2014-01-01) (2017), Definition: aspartokinase I / homoserine dehydrogenase I, Gene: thrA, thrA1, thrA2, Gene size: 2463 bp

Genome browser view showing gene structure, exons, introns, and annotations.

KEGGのパスウェイ

● 遺伝子間のつながりをグラフで表現したもの

- 代謝系
- 制御系

The image shows a screenshot of the KEGG Pathway Database website on the left and a detailed metabolic pathway diagram on the right. The website interface includes the KEGG logo, navigation tabs (KEGG2, PATHWAY, GENES, LIGAND, EXPRESSION, BRITE, XML, API, DBGE), and a list of categories such as Metabolism, Genetic Information Processing, Environmental Information Processing, Cellular Processes, and Human Diseases. The pathway diagram on the right is titled 'Glycolysis / Gluconeogenesis - Escherichia coli K-12 MG1655' and shows the conversion of glucose to pyruvate and back. Key enzymes are labeled with EC numbers, such as 2.7.1.1 (hexokinase), 1.1.1.1 (phosphoglycerate kinase), and 4.1.1.1 (phosphoenolpyruvate carboxylase). The diagram is divided into GLYCOLYSIS and GLUCONEOGENESIS sections.

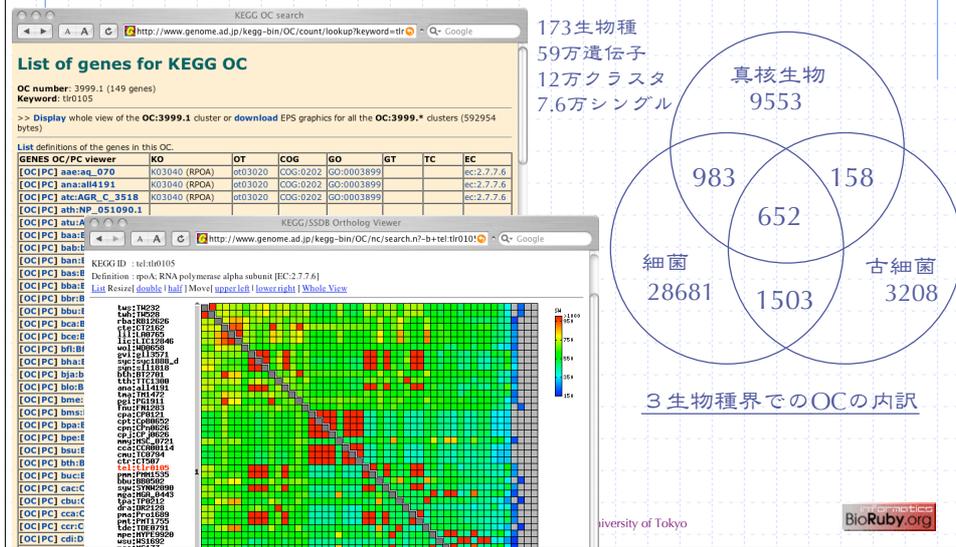
KO (KEGG Orthology)

● パスウェイへのアノテーションを効果的に進めるため、オーソログ対応の取れる遺伝子をキュレーション

- 酵素番号より細かい分類
- パスウェイ上の各ボックスの分類に相当
- 人手で行なうため時間がかかる

OCとSSDB

- 網羅的に候補を抽出するために全遺伝子ペアの配列類似性を計算しSSDBに格納、機械的にクラスタリング



KEGGを見てみよう

- <http://www.genome.jp/>
 - 遺伝子検索
 - GENESエントリ
 - LinkDB
 - パスウェイへのリンク
 - オーソログクラスタ
 - 遺伝子クラスタ
 - etc.

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.

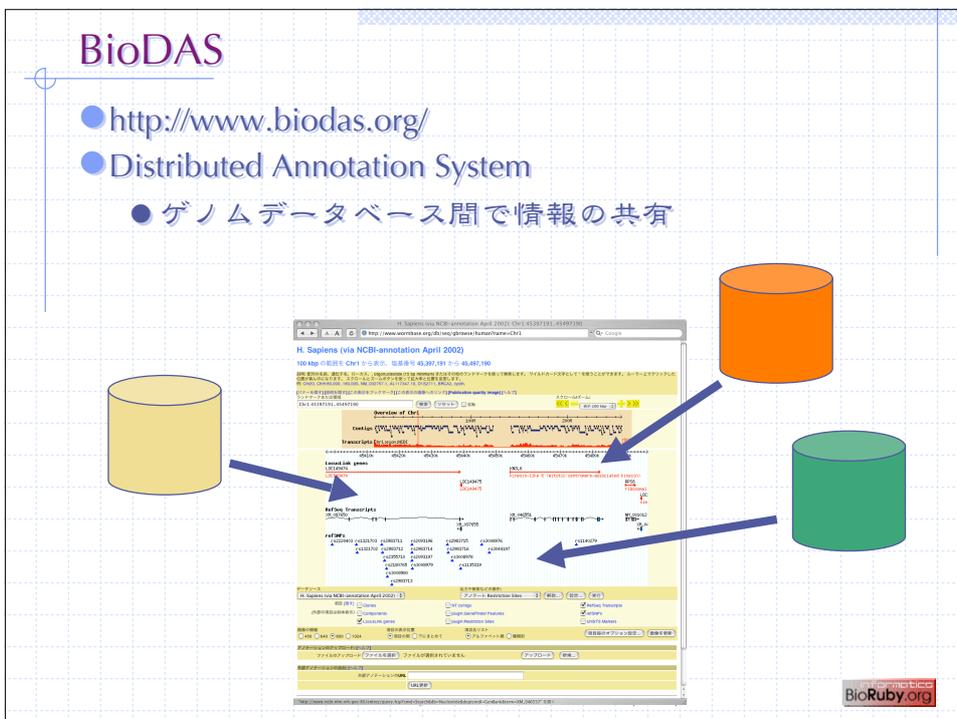
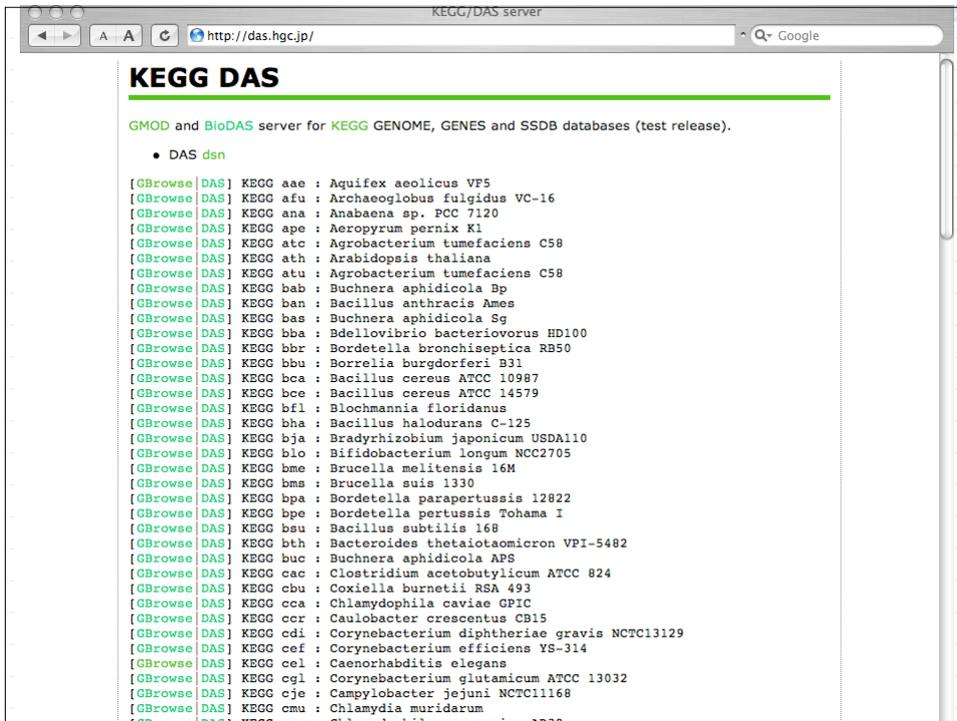
[Hide banner] [Hide instructions] [Bookmark this view] [Link to an image of this view] [Publication quality image] [Help]

Landmark or Region ana2364845..2384844 Flip **Scroll/Zoom:** <<< Show 20 kbp +>>>

Data Source KEGG ana - Anabaena sp. PCC 7120 **Dumps, Searches and other Operations:** Annotate Restriction Sites About... Configure... Go

KEGG DAS

- KEGGに含まれている生物種のゲノムアノテーション
 - 現在188生物種に対応
- GMODによるゲノムブラウザGBrowse
 - ユーザのアノテーションをアップロードできる
- DASによるXMLでのデータ取得
 - 真核生物だけでなく、バクテリア、古細菌にも網羅的に対応したDASサーバは少ない



GMOD

- Generic Model Organism Database
 - WormBase, FlyBase, SGD などのモデル生物で共通に利用可能でオープンなゲノムデータベースを開発
 - データベーススキーマ Chado (茶道)
 - ゲノムブラウザ GBrowse
 - アノテーションツール Apollo



KEGG DASを使ってみる

- <http://das.hgc.jp/>
 - 生物種を選択
 - クロモソームを選択
 - スクロール、ズーム
 - KEGGのアノテーションへリンク
 - ユーザのアノテーションをアップロード

DASのXMLデータ

```
<?xml version="1.0" standalone="yes"?>
<!DOCTYPE DASGFF SYSTEM "http://www.biodas.org/dtd/dasgff.dtd">
<DASGFF>
<GFF version="1.01" href="http://das.hgc.jp/cgi-bin/das/aae/features?segment=aae%3A10000%2C12000">
<SEGMENT id="aae" start="10000" stop="12000" version="1.0">
<FEATURE id="EC:1.1.1.22/50" label="1.1.1.22">
  <TYPE id="enzyme:KEGG" category="enzyme">enzyme:KEGG</TYPE>
  <METHOD id="enzyme">enzyme</METHOD>
  <START>11296</START>
  <END>12609</END>
  <SCORE></SCORE>
  <ORIENTATION></ORIENTATION>
  <PHASE>0</PHASE>
  <LINK href="http://www.genome.ad.jp/dbget-bin/w"
  <GROUP id="EC:1.1.1.22" type="EC" />
</FEATURE>
<FEATURE id="EC:2.6.1.11/46" label="2.6.1.11">
  <TYPE id="enzyme:KEGG" category="enzyme">enzyme:
  <METHOD id="enzyme">enzyme</METHOD>
  <START>10169</START>
  <END>11299</END>
  <SCORE></SCORE>
  <ORIENTATION></ORIENTATION>
  <PHASE>0</PHASE>
  <LINK href="http://www.genome.ad.jp/dbget-bin/w"
  <GROUP id="EC:2.6.1.11" type="EC" />
</FEATURE>
<FEATURE id="path:aae00040/51" label="aae00040">
  <TYPE id="pathway:KEGG" category="pathway">pathw
  <METHOD id="pathway">pathway</METHOD>
  <START>11296</START>
  <END>12609</END>
  <SCORE></SCORE>
  <ORIENTATION></ORIENTATION>
  <PHASE>0</PHASE>
  <LINK href="http://www.genome.ad.jp/dbget-bin/sh
  <GROUP id="path:aae00040" type="path" />
</FEATURE>
<FEATURE id="path:aae00040/48" label="aae00040">
  <TYPE id="pathway:KEGG" category="pathway">pathw
  <METHOD id="pathway">pathway</METHOD>
  <START>11296</START>
  <END>12609</END>
  <SCORE></SCORE>
  <ORIENTATION></ORIENTATION>
  <PHASE>0</PHASE>
  <LINK href="http://www.genome.ad.jp/dbget-bin/sh
  <GROUP id="path:aae00040" type="path" />
</FEATURE>
```

DASのURL

- 裏ではURLにより領域を指定、データ取得
 - XMLが返されるので必要な情報を抽出
- 塩基配列の取得
 - <http://das.hgc.jp/cgi-bin/das/eco/dna?segment=eco:101,200>
- アノテーション情報の取得
 - <http://das.hgc.jp/cgi-bin/das/sce/features?segment=l:5001,9000>
- BioRubyでデータ取得する方法は後ほど

KEGGの標準化

- KEGG DAS
 - GBrowse/DASによるゲノムアノテーション
 - ゲノム上での遺伝子の情報
- KEGG API
 - SOAP/WSDLによるKEGGへのアクセス
 - パスウェイ上の遺伝子群や類似遺伝子群の情報
- KGML
 - XMLによるKEGG PATHWAYの表現
 - 生物種ごとに遺伝子や化合物間のつながりの情報

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

BioRuby.org

KEGG API

- パスウェイやアノテーション情報をプログラムで取得
- カテゴリ (マニュアル参照)
 - メタ情報
 - DBGET
 - LinkDB
 - SSDB
 - Motif
 - KO, OC, PC
 - PATHWAY - get_genes_by_pathway(pathway_id)
 - GENES
 - GENOME
 - 具体的なメソッドや実際例は後ほど

© 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 for maps that are manually drawn and updated. KGML enables a facilities for computational analysis and modeling of protein n metabolic pathways contain two types of graph objects, how (compounds) are linked by a reaction in the KEGG pathway d pathways contain only the aspect of how boxes (proteins) are (KEGG Orthology) identifiers in the current KEGG system, bu are marked with EC numbers in the actual pathway diagram:

Documents

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

Data

- [KEGG reference metabolic pathways](#) (Last update: A)
- [KEGG reference regulatory pathways](#) (Last update: .)
- [KEGG metabolic pathways linked to KO](#) (Last update)
- [KEGG regulatory pathways linked to KO](#) (Last update)
- [KEGG organism-specific metabolic pathways](#)
Select organism:
- [KEGG organism-specific regulatory pathways](#)
Select organism:

```
<?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="path:map00010" 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="ec:1.2.1.3" type="enzyme" reaction="zn:r00710"
    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="ec:6.2.1.1" type="enzyme" reaction="zn:r00235"
    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="ec:1.2.1.5" type="enzyme" reaction="zn:r00711"
    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="path:map00650" type="map"
    link="http://www.genome.ad.jp/kegg/pathway/map/map00650.html">
    <graphics name="Butanoate metabolism" fgcolor="#000000" bgcolor="#
```

RubyとBioRuby

● Rubyとは

- 国産でフリーのオブジェクト指向スクリプト言語
- すっきりとした分かりやすい文法
 - 記号が少ない、便利なイテレータ(ブロック)
- 充実した組み込みクラスや標準添付ライブラリ
 - 数値、文字列、配列、ハッシュ、正規表現、...
 - クラスごとに使えるメソッドが豊富

● BioRubyとは

- Rubyのバイインフォマティクス用ライブラリ
- 2000年から主に日本で開発、KEGGなどもサポート

なぜRuby

- Rubyはどのような場面に向いているか
 - Perl同様、文字列処理や正規表現が得意
 - 遺伝子の配列やアノテーションなど文字情報が多い生物学のデータを扱うのに有利 → Perlの成功
- Rubyの代わりにPerl, Pythonでも全く問題はない
 - 普及率と使いやすさ（慣れ？）のトレードオフ
 - 使いたいライブラリがあるかどうか
 - JavaScript, PHP, SQLなども知っているると便利
- C, C++, Javaなどはコンパイルが必要
 - プログラムが長くなりがち → 手軽ではない
 - 高速な処理が可能、GUIアプリなどにも向く

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

BioRuby.org

Perlは？

- スクリプト言語の草分けでたぶん最も普及している
 - ユーザ数が多く、ライブラリも豊富
 - BioPerlがある
- でもデータ構造に弱い
 - 配列とハッシュの組み合わせ
 - Perl5以降のオブジェクト指向機能でカバー
 - 記号やオマジナイの増加 (リファレンス, bless, @ISA, ...)
- さらに基本的に使い捨てプログラム用の言語なので
 - 他人のPerlスクリプトは読めないことが多い
 - しばらく前に自分が書いたプログラムでも、

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

BioRuby.org

Open Bio* とは

- PerlにはBioPerl, PythonにはBiopython, JavaにはBioJava
言語ごとにライブラリ作成を進めているコミュニティが
- Open Bio Foundation
 - これら全体を支援し、BOSC, BioHackathonを開催

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

supernobio
BioRuby.org

BOSC

- ISMBと併催されるOpen Bio*のミーティング(BOF)
 - 進捗など成果発表を行ない、コミュニケーション



BioHackathon

- BOSCと異なりクローズド (企業スポンサー)
 - ホテルに缶詰になりハックしまくる

- 第1回 2002/01 アリゾナ
 - O'Reilly
- 第2回 2002/02 南アフリカ
 - Electric Genetics
- 第3回 2003/02 シンガポール
 - Apple Singapore

Schedule – hack hack hack

Starting		Ending	SUNDAY	MONDAY	TUESDAY	WEDNESDAY	THURSDAY	FRIDAY	SATURDAY
			24-Feb-02	25-Feb-02	26-Feb-02	27-Feb-02	28-Feb-02	1-Mar-02	2-Mar-02
7:00 AM	7:30 AM								
7:30 AM	8:00 AM								
8:00 AM	8:30 AM								
8:30 AM	9:00 AM			Hack Hack Hack	Hack Hack Hack	Hack Hack Hack	Hack Hack Hack	Hack Hack Hack	
9:00 AM	9:30 AM								
9:30 AM	10:00 AM								
10:00 AM	10:30 AM								
10:30 AM	11:00 AM			Coffee Break	Coffee Break	Coffee Break	Coffee Break	Coffee Break	
11:00 AM	11:30 AM			Hack Hack Hack	Hack Hack Hack	Hack Hack Hack	Hack Hack Hack	Hack Hack Hack	
11:30 AM	12:00 PM								
12:00 PM	12:30 PM								
12:30 PM	1:00 PM			Lunch	Lunch	Lunch	Lunch	Lunch	
1:00 PM	1:30 PM								
1:30 PM	2:30 PM			Hack Hack Hack	Off site Touring	Hack Hack Hack	Hack Hack Hack	Hack Hack Hack	
2:00 PM	2:30 PM				Dive / Snorkelling				
2:30 PM	3:00 PM				Kirstenbosch hike				
3:00 PM	3:30 PM				Waterfront				
3:30 PM	4:00 PM								
4:00 PM	4:30 PM			Coffee Break		Coffee Break			
4:30 PM	5:00 PM			Opening Reception & Hack Hack Hack	Hack Hack Hack	Hack Hack Hack	Township Tour & African Café	Coffee Break	
4:30 PM	5:00 PM			Cocktails & Snacks				Closing remarks	
5:00 PM	5:30 PM				Hack Hack Hack				
5:30 PM	6:00 PM								
6:00 PM	6:30 PM			Table Mountain					
6:30 PM	7:00 PM								
7:00 PM	7:30 PM			Dinner outside	Dinner at Hotel	Dinner at Hotel		Closing Party Outside	
7:30 PM	8:00 PM			at Castle					
8:00 PM	8:30 PM								
8:30 PM	9:00 PM				Hack Hack Hack	Hack Hack Hack	Hack Hack Hack		
9:00 PM	9:30 PM								
9:30 PM	10:00 PM								
10:00 PM	10:30 PM								
10:30 PM	11:00 PM								
11:00 PM	11:30 PM								
11:30 PM	12:00 AM								
12:00 AM	12:30 AM								
12:30 AM	1:00 AM								

Hacking Room



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

BioRuby.org

ラボ訪問やJunior BioHackathonも

- A*starの研究所のメンバーとディスカッション
- 14-16才くらいの子供たちへのレクチャー



BioRuby.org

スポンサー募集中...



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

BioRuby.org

BioRubyのカバーする領域

- Bio::Sequence, Bio::Location, Bio::Feature
 - 配列操作、スプライシング、ウィンドウサーチなど
- Bio::Registry, Bio::SQL, Bio::Fetch, Bio::FlatFile
 - OBDAによるデータベースへのアクセス
- Bio::DB
 - GenBank, SwissProt, KEGG などデータベースのパルザ
- Bio::Blast, Bio::Fasta (他 HMMER, MAFFT, Sosui, Psort etc.)
 - 解析アプリケーションの実行と結果のパルザ
- Bio::DAS, Bio::KEGG::API, Bio::DDBJ::XML
 - ウェブサービス
- Bio::Pathway, Bio::Relation
 - グラフ計算
- Bio::PubMed, Bio::Reference
 - 文献情報

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

BioRuby.org

配列操作

- スプライシング
- 相補鎖
- 翻訳
- 組成
- 分子量
- windowサーチ

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

require 'bio'

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

puts seq.to_fasta("foo", 60)
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

データベースへのアクセス

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

require 'bio'

serv = Bio::Registry.new
db = serv.get_database("swissprot")
entry = db.get_by_id("TETW_BUTFI")
puts entry
```

OBDA

KEGG API

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

require 'bio'

serv = Bio::KEGG::API.new
entry = serv.bget("sp:TETW_BUTFI")
puts entry
```

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

BioRuby.org

OBDA (Open Bio* Database Access)

- BioHackathon (2002, 2003)で制定
 - データベースのエントリを取得する仕組み
 - BioPerl, BioPython, BioJava, BioRubyで共通化
- 設定ファイル
 - ~/.bioinformatics/seqdatabase.ini
 - /etc/bioinformatics/seqdatabase.ini
 - <http://open-bio.org/registry/seqdatabase.ini>
 - チュートリアルを参照

© 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=biomysql
location=db.bioruby.org
dbname=biomysql
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



データベースエントリのパーサ

- 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
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

解析ソフトの実行

- 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.evalue, hsp.overlap
      end
    end
  end
end
```

BioRuby.org

サンプル集

- KEGG APIとKEGG DASをBioRubyから試してみる
 - SSDBを使ってホモログとモチーフを検索
 - KO, OC, PC で仲間の遺伝子を検索
 - データベースの情報やエントリの取得
 - パスウェイと遺伝子の対応
 - パスウェイに色をつける

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

BioRuby.org

Example1 : SSDBを使ってホモログとモチーフを検索

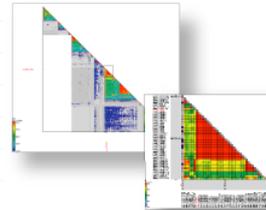
```
#!/usr/bin/env ruby

require 'bio'

# KEGG APIのサーバに接続
serv = Bio::KEGG::API.new

# 各生物種のbest-hit遺伝子を取得
homologs = serv.get_all_best_neighbors_by_gene("hsa:7368")

homologs.each do |hit|
  # hitした遺伝子名を取り出す
  gene = hit.genes_id2
  # KEGG APIで検索して遺伝子がモチーフを持っていれば
  if motifs = serv.get_motifs_by_gene(gene, "pfam")
    motifs.each do |motif|
      # 各モチーフのIDと説明を取り出してタブ区切りで表示
      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

Example2: KO, OC, PCで仲間の遺伝子

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

require 'bio'

serv = Bio::KEGG::API.new

# 遺伝子のKO番号のリストを取得
list = serv.get_ko_by_gene("eco:b0002")
# 例として1つ目のKO番号を使う
ko_id = list.first

# 同じKO (KEGG Orthology) がアサインされた遺伝子のリスト
ko_genes = serv.get_ko_members(ko_id)

# 同じOC (Ortholog Cluster) がアサインされた遺伝子のリスト
oc_genes = serv.get_all_oc_members_by_gene("hsa:7368")

# 同じPC (Paralog Cluster) がアサインされた遺伝子のリスト
pc_genes = serv.get_all_pc_members_by_gene("hsa:7368")

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

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

BioRuby.org

Example3: データベース情報やエントリ取得

```
#!/usr/bin/ruby

require 'bio'

serv = Bio::KEGG::API.new

# KEGGに現在含まれている生物種のリスト
orgs = serv.list_organisms
puts "# list of organisms in KEGG"
orgs.each do |entry|
  puts "#{entry.entry_id} #{entry.definition}"
end

# KEGGで利用可能なヒト(hsa)のパスウェイのリスト
list = serv.list_pathways("hsa")
puts "# list of pathways for human in KEGG"
list.each do |entry|
  puts "#{entry.entry_id} #{entry.definition}"
end

# 複数のKEGG GENESエントリを取得
puts "Human gene entries for gene 7368 and 7369"
puts serv.bget("hsa:7368 hsa:7369")
```

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

BioRuby.org

Example4: PATHWAYと遺伝子の対応

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

require 'bio'

serv = Bio::KEGG::API.new

# パスウェイ上の遺伝子のリスト
genes = serv.get_genes_by_pathway("path:hsa00020")
puts "# genes on human's pathway 00020"
genes.each do |gene|
  puts gene
end

# EC番号を遺伝子名に対応づけ
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

Example5: PATHWAYに色を付ける

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

require 'bio'

serv = Bio::KEGG::API.new

# パスウェイに赤枠でマークをつける
objs = ['eco:b0002', 'cpd:C00263']
url1 = serv.mark_pathway_by_objects('path:eco00260', objs)
puts "# URL for marked pathway:", url1

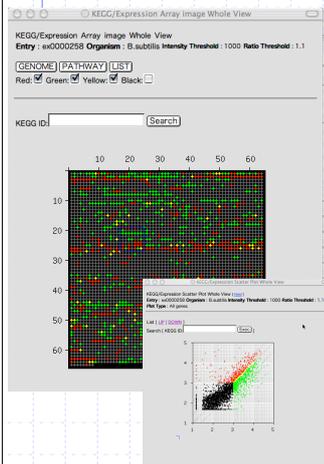
# パスウェイに色をつける
fg_list = ['blue', 'green']
bg_list = ['#ff0000', 'yellow']
url2 = serv.color_pathway_by_objects('path:eco00260',
  objs, fg_list, bg_list)
puts "# URL for colored pathway:", url2

# 結果の画像を保存する
serv.save_image(url1, "example5-1.gif")
serv.save_image(url2, "example5-2.gif")
```

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

BioRuby.org

Example6: 遺伝子発現をパスウェイにマッピング

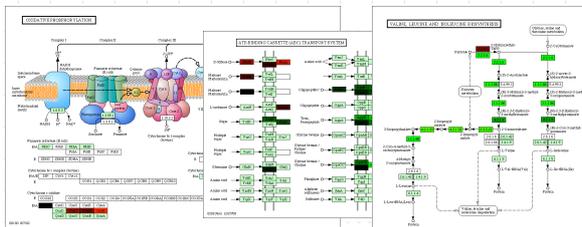


```
serv = Bio::KEGG::API.new

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) # 遺伝子名と色の対応
end

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



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

biokuby.org

Example7: PDB へのマッピング

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

require 'bio'

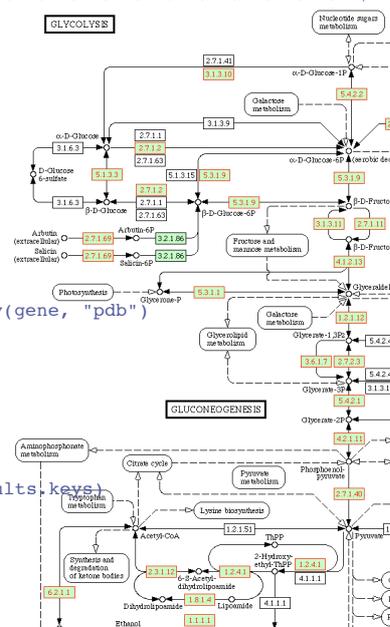
serv = Bio::KEGG::API.new

# 指定したいパスウェイ上の遺伝子のリスト
path = ARGV.shift || "path:eco00010"
genes = serv.get_genes_by_pathway(path)

# PDBにリンクのある遺伝子を検索
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

# 色付き画像を生成
url = serv.mark_pathway_by_objects(path, results.keys)

# 画像を保存
serv.save_image(url, "example7.gif")
```



© 2004 Toshiaki Katayama, Human Genome Center

Example8: DASからの情報取得

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

require 'bio'

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

# 酵母のI番染色体1001塩基から2000塩基の範囲を指定、DNA配列を取得
segment = Bio::DAS::SEGMENT.region("I", 1001, 2000)
list = serv.get_dna("sce", segment)
list.each do |dna|
  puts dna.sequence
end

# 上記と同じ範囲で遺伝子アノテーションを取得
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

BioRuby.org

なぜオブジェクト指向か

- そんなに大げさなものではなく、単に便利だから
 - BioPerl, Biopython, BioJava, BioRuby 全て...
- オブジェクトはデータの入れもの (構造体)
 - さらにデータの操作方法 (メソッド) を定義
 - 化合物オブジェクトには名前、構造式、機能説明などのデータ (情報) が入っている
 - 組成や分子量を計算するメソッドなど
- オブジェクトに共通の特徴を括り出したものがクラス
 - 抽象化 (モノをどう捉えるか、理解するか)
 - 複数のクラスに共通の特徴を括り出したらスーパークラス (これを各クラスで継承し差分を書く)

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

BioRuby.org

クラスを作ってみる - Ruby入門

- いれものであるCpdクラスを宣言
 - 名前と構造式と説明を覚えておく

```
class Cpd
  attr_accessor :name, :definition, :formula
end
```

```
a = Cpd.new
a.name = "ATP"
a.definition = "Adenosine 5'-triphosphate"
a.formula = "C10H16N5O13P3"

puts a.name      # => "ATP"
puts a.definition # => "Adenosine 5'-tri.."
puts a.formula   # => "C10H16N5O13P3"
```

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

BioRuby.org

クラスを作ってみる - Ruby入門

- 本当なら...

```
class Cpd
  def initialize
    @formula = nil
  end

  def formula=(val)
    @formula = val
  end

  def formula
    return @formula
  end
end
```

```
cpd = Cpd.new
cpd.formula = "C10H16N5O13P3"
puts cpd.formula
# => "C10H16N5O13P3"
```

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

BioRuby.org

クラスを作ってみる - メソッド定義

- 構造式から組成を抽出するメソッドを追加

```
class Cpd

  attr_accessor :name, :definition, :formula

  # 組成を覚えておくHash型の変数を準備
  def initialize
    @comp = Hash.new
  end

  # 組成を抽出するメソッド "C10H16N5O13P3"
  def composition
    @formula.scan(/([A-Z]+)(\d+)/) do |a, b|
      @comp[a] = b.to_i
    end
    return @comp
  end

end
```

{ "C" => 10, "H" => 16,
"N" => 5, "O" => 13,
"P" => 3 }

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

BioRuby.org

クラスを作ってみる - もう少し機能追加

- 組成から全体の分子量を計算するメソッドを追加

```
class Cpd

  # 手抜きのために分子量の表も内蔵
  MW = { 'C'=>12.011, 'H'=>1.00794, 'N'=>14.00674,
         'O' => 15.9994, 'P' => 30.973762 }

  def initialize
    @comp = Hash.new
  end

  attr_accessor :name, :definition, :formula

  def composition
    @formula.scan(/([A-Z]+)(\d+)/) do |a, b|
      @comp[a] = b.to_i
    end
    return @comp
  end

  def mol_weight
    total = 0.0
    composition.each do |elem, i|
      total += MW[elem] * i
    end
    return total
  end

end
```

BioRuby.org

クラスを作ってみる - 実行

- 以下の部分をつけてcpdtest1.rbに保存し、実行

```
# (続き)
cpd = Cpd.new
cpd.name = "ATP"
cpd.definition = "Adenosine 5'-triphosphate"
cpd.formula = "C10H16N5O13P3"
```

```
puts cpd.name
puts cpd.definition
puts cpd.formula
puts cpd.mol_weight
```

```
% ruby cpdtest1.rb
ATP
Adenosine 5'-triphosphate
C10H16N5O13P3
507.184226
```

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

BioRuby.org

クラスを作ってみる - ライブラリ化

- Cpdクラスの定義部分だけをcpd.rbに括り出し再利用

```
require 'cpd'

cpd = Cpd.new
cpd.name = "ATP"
cpd.definition = "Adenosine 5'-triphosphate"
cpd.formula = "C10H16N5O13P3"
```

```
puts cpd.name
puts cpd.definition
puts cpd.formula
puts cpd.mol_weight
```

```
% ruby cpdtest2.rb
ATP
Adenosine 5'-triphosphate
C10H16N5O13P3
507.184226
```

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

BioRuby.org

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

BioRuby.org