東京工業大学 大学院情報理工学研究科 角田将典、石田貴士、秋山泰

2015年3月20日

並列配列相同性検索プログラム 「GHOST-MP」講習会(実習編)





SCLS計算機システム講習会









角田 将典 かくた まさのり 石田 貴士 いしだ たかし 秋山 泰 あきやま ゆたか

東京工業大学 大学院情報理工学研究科 計算工学専攻





- GHOST-MPの利用方法の習得
 - GHOST-MPの基本的な利用方法
 - SCLS計算機システム上での利用
 - スケジューラを介した実行方法
- GHOST-MPによるメタゲノムデータに対する 相同性検索結果の解析
 - 分類群に基づく解析
 - 遺伝子オーソログに基づく解析

実習前の確認



- SCLS計算機システムにログインできること
- ログインノードにおける基本的なファイル操作が できること
- ウェブブラウザでAdobe Flash Playerが
 利用可能なこと
 - * GHOST-MPの実行とは直接関係はないが、 GHOST-MPの検索結果の可視化の際に必要

アジェンダ



- GHOST-MP実行のながれ
- STEP 1: 実習の準備
 - STEP 1-1: SCLS計算機システムへのログイン
 - STEP 1-2: 必要なファイルのコピーと展開
- STEP 2: ジョブスケジューラの利用方法
- STEP 3: GHOST-MPによる配列相同性検索
 - STEP 3-1: GHOST-MPのコンパイル
 - STEP 3-2: データベースのインデックス作成(ghostmp_makedb)
 - STEP 3-3: 配列相同性検索(ghostmp_search)
- STEP 4: メタゲノムデータに対する相同性検索結果の解析
 - STEP 4-1: 検索結果の集計
 - STEP 4-2: 分類群に基づく解析
 - STEP 4-3: 遺伝子オーソログに基づく解析

本日の実習内容











GHOST-MP

BLAST のように遠縁のホモログを検出可能な配列相同性検索ツールである GHOSTXアルゴリズムを並列化したもの

GHOST-MP

- MPIとOpenMPによるハイブリッド並列化 MPIによる計算ノード間の並列化、OpenMPによる計算ノード内 の並列化を行い、クエリ配列のデータ並列処理を行う
- ファイル I/O
 MPI-IOを利用し、並列入出力を行う
- 事前にデータベースのインデックス作成が必要 GHOSTXアルゴリズムは、高速な検索を実現するために、接尾辞 配列(Suffix Array)というデータ構造を事前に作成しておく



- GHOST-MP関連
 - GHOSTX: Suzuki et al., (2014) *PLoS ONE* 9(8):e103833
 - 接尾辞配列を用いたアラインメント候補位置の高速探索
 - GHOST-MP: Kakuta et al., (in preparation)
 - GHOSTXの分散メモリ環境版、
- 当グループの他の配列相同性検索関連
 - GHOXTM: Suzuki et al., (2012) PLoS ONE 7(5): e36060
 - GPUを用いた相同配列検索
 - GHOSTZ: Suzuki et al., (in press) doi: 10.1093/bioinformatics/btu780
 - 部分文字列のクラスタリングによるアラインメント候補位置の高速探索
 - GHOSTZ-GPU: Suzuki et al., (in preparation)
 - GHOSTZのGPU版







GHOST-MP

MPI-IOによる並列ファイル入出力



MPIライブラリによる効率的なファイルアクセス

実習のながれ



実習用にあらかじめ用意した クエリ配列 と データベース配列 を使って、 配列相同性検索の並列分散処理を実行し、検索結果の解析を行う。



GHOST-MP コマンドの実行のながれ





アジェンダ



- GHOST-MP実行のながれ
- STEP 1: 実習の準備
 - STEP 1-1: SCLS計算機システムへのログイン
 - STEP 1-2: 必要なファイルのコピーと展開
- STEP 2: ジョブスケジューラの利用方法
- STEP 3: GHOST-MPによる配列相同性検索
 - STEP 3-1: GHOST-MPのコンパイル
 - STEP 3-2: データベースのインデックス作成(ghostmp_makedb)
 - STEP 3-3: 配列相同性検索(ghostmp_search)
- STEP 4: メタゲノムデータに対する相同性検索結果の解析
 - STEP 4-1: 検索結果の集計
 - STEP 4-2: 分類群に基づく解析
 - STEP 4-3: 遺伝子オーソログに基づく解析



• sshによるログイン

- ログインノード: hpci-scls.riken.jp
- 事前に用意したsshの秘密鍵を使って、 仮想端末などからログインする
 - Windows: PuTTYやTera Term
 - LinuxやMac OS X: sshコマンド



\$ ssh username@hpci.scls.riken.jp SCLS System Info: Service is available. Last login: Thu Mar 10 14:34:45 2015 from foo.ac.jp		4
# SCLS System Information Date : Mar. 9 #	# , 2015 #	+ + +
	#	#
# Welcome to SCLS System	#	‡
#	#	ŧ
# If you have any questions or need for further assistance, please send # measure to gala admenitor in	da #	‡ +
# message to sets-admeriken.jp	#	+ +
<pre># For SCLS information, please refer to https://hpci-scls.riken.jp/</pre>	#	‡
#	#	+
#	#	+
[username@scls ~]\$		

* 仮想端末の表示は環境設定等のため、必ずしも一致しない

アジェンダ



- GHOST-MP実行のながれ
- STEP 1: 実習の準備
 - STEP 1-1: SCLS計算機システムへのログイン
 - STEP 1-2: 必要なファイルのコピーと展開
- STEP 2: ジョブスケジューラの利用方法
- STEP 3: GHOST-MPによる配列相同性検索
 - STEP 3-1: GHOST-MPのコンパイル
 - STEP 3-2: データベースのインデックス作成(ghostmp_makedb)
 - STEP 3-3: 配列相同性検索(ghostmp_search)
- STEP 4: メタゲノムデータに対する相同性検索結果の解析
 - STEP 4-1: 検索結果の集計
 - STEP 4-2: 分類群に基づく解析
 - STEP 4-3: 遺伝子オーソログに基づく解析

実習の準備



実習に必要なファイルのコピーと展開

- SCLS計算機システム上の 「/home/kakuta/ghostmp_koshukai.tar.gz」をコピーする
- 必要ならば適宜ディレクトリを作成して、そのディレクトリの中で作業する





- data/:
 実習に用いるメタゲノムデータ
- ghostmp-1.3.3/: GHOST-MPのソースコード
- precomp/: 事前に計算されたGHOST-MPの結果(時間がかかるため)
- sample/: 実習で作成するスクリプトなどのサンプル
- script/: GHOST-MPの出力結果の解析に用いるスクリプト
- workspace/:
 このディレクトリで実習を行うことを想定



- GHOST-MP実行のながれ
- STEP 1: 実習の準備
 - STEP 1-1: SCLS計算機システムへのログイン
 - STEP 1-2: 必要なファイルのコピーと展開
- STEP 2: ジョブスケジューラの利用方法
- STEP 3: GHOST-MPによる配列相同性検索
 - STEP 3-1: GHOST-MPのコンパイル
 - STEP 3-2: データベースのインデックス作成(ghostmp_makedb)
 - STEP 3-3: 配列相同性検索(ghostmp_search)
- STEP 4: メタゲノムデータに対する相同性検索結果の解析
 - STEP 4-1: 検索結果の集計
 - STEP 4-2: 分類群に基づく解析
 - STEP 4-3: 遺伝子オーソログに基づく解析





ジョブの実行方法(1)



- SCLS計算機システムでのジョブの実行には、キューイングシステムを利用する
- 次のように実行するジョブを記述したスクリプトを作成し、 pjsubコマンドでジョブのサブミットを行う

[username@scls workspace]\$ pwd /home/username/ghostmp_koushukai/workspace [username@scls workspace]\$ mkdir job_submit [username@scls workspace]\$ cd job_submit [username@scls job_submit]\$ vim run.sh

スクリプトの例

run.sh

#!/bin/sh
#PJM -L "rscgrp=small"
#PJM -L "node=1"
#PJM -L "elapse=00:10:00"
sleep 120
echo "test message"





キューイングシステムの主なコマンド

- pjsub: ジョブの投入
 - pjsub sample_job.sh
- pjstat:投入されたジョブの状態表示
 - pjstat
- pjdel:投入されたジョブの削除
 - pjdel job_id



run.shの投入例

[username@scls]\$ pjsub run.sh 「TWE0」 PTM 00000 - i.e.h. Teh. 10047h. はしし	
[INFO] PJM 0000 pjsub Job 19247 submitted.	
[username@scls]\$ pjstat 投入したジョブの状態を出力する	
ACCEPT QUEUED STGIN READY RUNING RUNOUT STGOUT HOLD ERROR TOTAL	
0 0 0 0 1 0 0 0 1	
s 0 0 0 0 1 0 0 0 0 1	
JOB_ID JOB_NAME MD ST USER START_DATE ELAPSE_LIM NODE_REQUIRE	
19247 run.sh NM RUN username 03/11 11:42:58 0000:05:00 1	
[usename@scls]\$ pjstat	
ジョブが終了するとpjsubで表示されなくなる	
ACCEPT QUEUED STGIN READY RUNING RUNOUT STGOUT HOLD ERROR TOTAL	
0 0 0 0 0 0 0 0 0 0	
s 0 0 0 0 0 0 0 0 0 0 0	
[username@scls]\$ ls 🐂 「「」」)	ל ל
run.sh run.sh.e19247 run.sh.o19247 がrun.sh.o <job_id>とrun.sh.e<job_id>に出力され</job_id></job_id>	れる
[username@scls]\$ cat run.sh.o19247	
test message	
[username@scls]\$	



[userna	me@scls]	\$ pjstat								
ACCEP	T QUEUED	STGIN	READY	RUNING	RUNOUT	STGOUT	HOLD	ERROR	TOTAL	r + 1
	0 0	0	0	1	0	0	0	0	1	[I]
S	0 0	0	0	1	0	0	0	0	1	[2]
JOB_ID 19247	JOB_ run.	NAME MI sh NN	OSTU MRUNU	JSER Isername	START_ e 03/11	_DATE 11:42:58	ELAP 0000	SE_LIM : :05:00	NODE_REQUIRE 1	[3]

[1] ある状態のジョブ数
[2] サブジョブを考慮したジョブ数 (本実習ではサブジョブは使用しない)
[3] 各ジョブの情報

[1]および[2]	内容
ACCEPT	受入中のジョブ数
QUEUED	ステージイン待機中のジョブ数
STGIN	ステージイン中のジョブ数
READY	実行待機中のジョブ数
RUNING	実行中のジョブ数
RUNOUT	ジョブ終了中のジョブ数
STGOUT	ステージアウト中のジョブ数
HOLD	ユーザー指示で待機中のジョブ数
ERROR	エラーのため待機中のジョブ数
TOTAL	全ジョブ数



[ι	ısername	@scls]\$	pjstat								
	ACCEPT	QUEUED	STGIN	READY	RUNING	RUNOUT	STGOUT	HOLD	ERROR	TOTAL	r 1 1
	0	0	0	0	1	0	0	0	0	1	[]
Ŋ	0	0	0	0	1	0	0	0	0	1	[2]
JC)B_ID	JOB_N	AME M	D ST (JSER	START	_DATE	ELAP	SE_LIM 1	NODE_REQUIR	E
19	9247	run.s	h N	M RUN ı	username	e 03/11	11:42:58	3 0000	:05:00	1	[3]

[3]	内容	ST	内容
JOB_ID	キューイングシステムによってジョブに	ACC	受入中
	割り振られたID	RJT	棄却中
JOB_NAME	ジョブの名前。ユーザーが指定できる。 未指定の場合は、スクリプト名。	QUE	実行待機中
MD	ジョブの種類(通堂のジョブ ステップジョブ	SIN	ステージイン中
	バルクジョブ)	RDY	実行待機中(SIN後)
ST	ジョブの状態	RNA	計算資源確保中
USER	ジョブを投入したユーザー	RUN	実行中
START_DATE	ジョブの実行開始日時	RNO	ステージアウト待機中
ELAPSE_LIM	ジョブがRUN状態になってからのの経過時間	SOT	ステージアウト中
NODE_REQUIRE	要求ノード	EXT	ジョブ終了処理完了
		CCL	ジョブ実行中止
		HLD	ユーザーの指示で待機中
•		ERR	エラー状態



- ・ 誤ったジョブを投入してしまったときなどは、
 ジョブの削除が可能
 - pjdel:投入されたジョブの削除
 - pjdel job_id
 - 削除したいジョブのJOB_IDをpjstatなどで調べて、
 pjdelコマンドの引数として指定する
 - 待機中のジョブも、実行中のジョブの削除可能

	_										
luserr	username@scls]\$ pistat										
Laberr	.o.the	SDCTDI	PJDeac								
	'DT	OTETED	STGIN	READY	RINTNG	RINOIT	STGOIT	нот.р	EBBOB		
IICCI			DIGIN		1011110	10011001	010001		BIUROIC	IOIAH	
	0	0	0	0	1	0	0	0	0	1	
	~	0	0		1	0	0	0	0	1	
S	0	U	0	0		U	0	0	0	\perp	
JOB_II		JOB_NZ	AME MI) ST U	JSER	START_	DATE	ELAI	SE_LIM	NODE_REQUIRE	
10017					100000	02/11	11.10.5	0 0000		1	
19241		run.si	.1 1019	I RON U	isername	2 03711	11.47.0	8 0000	1.02.00	L	
luserr	ame		nidel 1	9247							
LUDCII	lance	SPCTD]Å	Place 1	- 7 2 1 7							



- GHOST-MP実行のながれ
- STEP 1: 実習の準備
 - STEP 1-1: SCLS計算機システムへのログイン
 - STEP 1-2: 必要なファイルのコピーと展開
- STEP 2: ジョブスケジューラの利用方法
- STEP 3: GHOST-MPによる配列相同性検索
 - STEP 3-1: GHOST-MPのコンパイル
 - STEP 3-2: データベースのインデックス作成(ghostmp_makedb)
 - STEP 3-3: 配列相同性検索(ghostmp_search)
- STEP 4: メタゲノムデータに対する相同性検索結果の解析
 - STEP 4-1: 検索結果の集計
 - STEP 4-2: 分類群に基づく解析
 - STEP 4-3: 遺伝子オーソログに基づく解析

アジェンダ



- STEP 3: GHOST-MPによる配列相同性検索
 - STEP 3-1: GHOST-MPのコンパイル
 - STEP 3-2: データベースのインデックス作成(ghostmp_makedb)
 - STEP 3-3: 配列相同性検索(ghostmp_search)

STEP 3-3までの作業で得られる結果 各クエリ配列の相同配列とそのアラインメント情報

query_name sbj_name 0.714286 14 2 1 32 73 105 116 2.88961 29.4819

各カラムが示すもの

#query_name sbj_name identity(%) alignment_length num_mismatch num_gap query_start query_end sbj_start sbj_end e_value bit_score

コンパイル



Requirement

- OpenMPに対応したコンパイラ
- MPI Library
- tr1またはBoost C++ Library

ghostmp_makedb

ghostmp_search



\$ ls ghostmp_makedb ghostmp_search
ghostmp_makedb

コンパイルに成功すると ghostmp_makedbと ghostmp_search が作成される



- GHOST-MP実行のながれ
- STEP 1: 実習の準備
 - STEP 1-1: SCLS計算機システムへのログイン
 - STEP 1-2: 必要なファイルのコピーと展開
- STEP 2: ジョブスケジューラの利用方法
- STEP 3: GHOST-MPによる配列相同性検索
 - STEP 3-1: GHOST-MPのコンパイル
 - STEP 3-2: データベースのインデックス作成(ghostmp_makedb)
 - STEP 3-3: 配列相同性検索(ghostmp_search)
- STEP 4: メタゲノムデータに対する相同性検索結果の解析
 - STEP 4-1: 検索結果の集計
 - STEP 4-2: 分類群に基づく解析
 - STEP 4-3: 遺伝子オーソログに基づく解析



• 注意

- 計算時間短縮のため、「データベースのインデックス作成」と「配列相同性検索」では、実際に解析に用いられるファイルの一部が用いられています。
- 「検索結果の集計」からは、事前計算された結果を利用して頂きますので、ファイルを取り違えないように注意して下さい。

インデックスの作成 [ghostmp_makedb] (入力ファイル)





配列相同性検索の参照先となる FASTA フォーマットのアミノ酸配列の データベースファイルからインデックスを作成する

FASTA フォーマットとは

>gi|66816243|ref|XP_642131.1| hypothetical protein (*>"で始まるヘッダ行 配列の説明 1行で書く MASTQNIVEEVQKMLDTYDTNKDGEITKAEAVEYFKGKKAFNPERSAIYLFQVYDKDNDGKITIKELAGDIDFDKALKEY KEKQAKSKQQEAEVEEDIEAFILRHNKDDNTDITKDELIQGFKETGAKDPEKSANFILTEMDTNKDGTITVKELRVYYQK VQKLLNPDQ 配列データ アミノ酸を一文字表記で表している 改行 OK

データベースサイズが大きい場合(数GB~数+GB) #個 のチャンクに分割して扱う。 それぞれのチャンクについて 5つのファイル (.ind .inf .nam .off .seq) が作成される。

インデックスの作成 [ghostmp_makedb] (実行準備)



ghostmp_makedb

ghostmp_makedb -i INPUT -o OUTPUT - INPUT: データベース配列 - OUTUPT: 作成するインデックスのファイル名

ジョブスクリプト (逐次ジョブ)の作成

[username@scls job_submit]\$ pwd /home/kakuta/ghostmp_koushukai/workspace/job_submit [username@scls job_submit]\$ mkdir ../makedb && cd \$_ [username@scls makedb]\$ vim makedb.sh

makedb.sh





ジョブの投入と 状態確認

[user [INFC [user	rname)] PJ rname	@scls ma M 0000 j @scls ma	akedb]\$ pjsub Jc akedb]\$	pjsub b 1056 pjstat	makedb. 9 submi	sh tted.				
ACC	EPT	QUEUED	STGIN	READY	RUNING	RUNOUT	STGOUT	HOLD	ERROR	TOTAL
	0	0	0	0	1	0	0	0	0	1
S	0	0	0	0	1	0	0	0	0	1
JOB_I	D	JOB_N	AME MD) ST U	ISER	START_	_DATE	ELAP	SE_LIM	NODE_REQUIRE
10569)	maked	b.sh NM	I RUN u	Isername	e 10/24	14:49:49	5 0000	:10:00	1



[username@scls makedb]\$ cat makedb.sh.	. 0*	
The number of chunks :1 Max length of a chunk : 50476865 Total database length : 50369581 Total number of sequences : 107283	インデックスのチャンク数1 チャンクの最大長 配列の合計長 全配列数	

インデックスの作成 [ghostmp_makedb] (出力ファイルの確認)







- GHOST-MP実行のながれ
- STEP 1: 実習の準備
 - STEP 1-1: SCLS計算機システムへのログイン
 - STEP 1-2: 必要なファイルのコピーと展開
- STEP 2: ジョブスケジューラの利用方法
- STEP 3: GHOST-MPによる配列相同性検索
 - STEP 3-1: GHOST-MPのコンパイル
 - STEP 3-2: データベースのインデックス作成(ghostmp_makedb)
 - STEP 3-3: 配列相同性検索(ghostmp_search)
- STEP 4: メタゲノムデータに対する相同性検索結果の解析
 - STEP 4-1: 検索結果の集計
 - STEP 4-2: 分類群に基づく解析
 - STEP 4-3: 遺伝子オーソログに基づく解析

配列相同性検索 [ghostmp_search] (入力ファイル)





※ クエリファイルの配列は、次世代シーケンサーで得られた HMPの頬粘膜のDNA塩基配列データ。



ghostmp_search

ghostmp_search -a NUM_THREADS -q QUERY_TYPE -i INPUT -d DB -o OUTPUT

- NUM_THREADS: スレッド数
- QUERY_TYPE: クエリが塩基配列ならば"d"、アミノ酸配列ならば"p"を指定する
- INPUT: クエリ配列
- DB: ghostmp_makedbで作成したインデックスのprefix
- OUTUPT: 検索結果の出力先

ジョブスクリプト(MPI並列ジョブ)

[username@scls workspace]\$ mkdir search && cd \$_
[username@scls search]\$ vim search.sh

search.sh



配列相同性検索 [ghostmp_search] (ジョブの投入)





ジョブ終了確認

[ι	[username@scls search]\$ pjstat										
	ACCEPT Ç	QUEUED	STGIN	READY	RUNING	RUNOUT	STGOUT	HOLD	ERROR	TOTAL	
	0	0	0	0	0	0	0	0	0	0	
S	0	0	0	0	0	0	0	0	0	0	

配列相同性検索 [ghostmp_search] (出力ファイルの確認)





結果確認1 – 出力フ	アイル					
[username@scls search]	\$ ls -1					
total 2880						
-rw-rr 1 username g	group 44553	77 Oct	24	16:24	search_result	ghostmp_searchの検索結果
-rw-rr 1 username g	group 3	24 Oct	24	16:20	search.sh	
-rw-rr 1 username g	group	0 Oct	24	16:24	search.sh.e10583	
-rw-rr 1 username g	group	0 Oct	24	16:24	search.sh.o10583	ghostmp_searchの標準出力

配列相同性検索 [ghostmp_search] (出力内容の確認)





※検索結果は、1行ごとにTAB区切りで出力される

クエリ配列名、ヒット配列名、一致率、アラインメント長、不一致数、ギャップ数、 クエリ配列におけるアラインメント開始位置、その終了位置、 ヒット配列におけるアラインメント開始位置、その終了位置、 E-value、bitスコア

「京」での実行スクリプト例



「京」とSCLS計算機システムは互換性の高いシステムだが、 「京」ではジョブに必要なファイルを転送するための記述が必要



#!/bin/sh #PJML "rscgrp=small" 「京」での実行スクリプト例	
#PJML "elapse=0:10:00"	
#PJML "node=2"	
#PJMmpi "use-rankdir"	
<pre>#PJMstgin "rank=*//ghostmp-1.3.3/src/ghostmp_search %r:./ghostmp_search"</pre>	
<pre>#PJMstgin "rank=0//data/query 0:/"</pre>	
<pre>#PJMstgin "rank=0/makedb/out/db* 0:/db/"</pre>	
<pre>#PJMstgout "rank=0 0:/search_result ./search_result"</pre>	

ステージインとステージアウトに関する記述

アジェンダ



- GHOST-MP実行のながれ
- STEP 1: 実習の準備
 - STEP 1-1: SCLS計算機システムへのログイン
 - STEP 1-2: 必要なファイルのコピーと展開
- STEP 2: ジョブスケジューラの利用方法
- STEP 3: GHOST-MPによる配列相同性検索
 - STEP 3-1: GHOST-MPのコンパイル
 - STEP 3-2: データベースのインデックス作成(ghostmp_makedb)
 - STEP 3-3: 配列相同性検索(ghostmp_search)
- STEP 4: メタゲノムデータに対する相同性検索結果の解析
 - STEP 4-1: 検索結果の集計
 - STEP 4-2: 分類群に基づく解析
 - STEP 4-3: 遺伝子オーソログに基づく解析

アジェンダ



- STEP 4: メタゲノムデータに対する相同性検索結果の解析
 - STEP 4-1: 検索結果の集計
 - STEP 4-2: 分類群に基づく解析
 - STEP 4-3: 遺伝子オーソログに基づく解析

STEP 4-3までの作業で得られる結果



STEP4-2: サンプル内のphylumの割合



STEP4-3: サンプル内の遺伝子が 関わる代謝パスウェイの可視化







result.buccal_mucosa: ヒトの頬粘膜(HMPのSRS011090)の検索結果 result.supragingival_plaque: ヒトの歯肉縁上の歯垢(HMPのSRS044366)の検索結果 (HMP: Human Microbiome Project。ヒトの体表・体内のメタゲノム研究プロジェクト)



supragingival_plaqueでも同様に解析を実行する





配列相同検索の結果

summarize_search_result.py

出力ファイル

解析を実行する 標準出力

[user1@scls summary]\$ python ../../script/summarize_search_result.py ¥
../../precomp/result.buccal_mucosa

2013/11/21	14:44:21	START	Start KEGG Analyzer
2013/11/21	14:44:21	START	Loading Gi-Taxid map file
2013/11/21	14:45:27	END	54380376 genes loaded.
2013/11/21	14:45:27	START	Loading Taxonomy root map file
2013/11/21	14:45:33	END	919194 species loaded.
2013/11/21	14:45:33	START	Loading KO-Enzyme map file
2013/11/21	14:45:33	END	4808 KOs loaded.
2013/11/21	14:45:33	START	Loading USCG map file
2013/11/21	14:45:33	END	36 USCGs loaded.
2013/11/21	14:45:33	START	Loading KEGG genes file
2013/11/21	14:46:22	END	8782317 genes loaded.
2013/11/21	14:46:24	START	Loading Blast result
2013/11/21	14:46:25	END	15000 blast results loaded.
2013/11/21	14:46:25	START	Normalizing
2013/11/21	14:46:25	START	1. count genes
2013/11/21	14:46:26	END	done.
2013/11/21	14:46:26	START	2. count USCGs
2013/11/21	14:46:26	END	done.
2013/11/21	14:46:26	START	3. normalizing
2013/11/21	14:46:28	END	done.
2013/11/21	14:46:28	END	Normalize.
2013/11/21	14:46:28	END	End KEGG Analyzer





配列相同検索の結果

summarize_search_result.py



出力ファイル(主なもの)

ファイル名	内容
<search_result>.genes_freq</search_result>	遺伝子の出現頻度
<search_result>.ko_ratio</search_result>	各KO (KEGG Orthology)の割合
<search_result>.phylum_ratio</search_result>	phylum の割合

集計方法

- 検索結果全体から、各遺伝子の出現頻度を推定する。
- 各遺伝子に対するKEGGのアノテーションからKOについて集計する。
- pylumとgenusについてはマーカー遺伝子のみを用いて推定を行う。

アジェンダ



- GHOST-MP実行のながれ
- STEP 1: 実習の準備
 - STEP 1-1: SCLS計算機システムへのログイン
 - STEP 1-2: 必要なファイルのコピーと展開
- STEP 2: ジョブスケジューラの利用方法
- STEP 3: GHOST-MPによる配列相同性検索
 - STEP 3-1: GHOST-MPのコンパイル
 - STEP 3-2: データベースのインデックス作成(ghostmp_makedb)
 - STEP 3-3: 配列相同性検索(ghostmp_search)
- STEP 4: メタゲノムデータに対する相同性検索結果の解析
 - STEP 4-1: 検索結果の集計
 - STEP 4-2: 分類群に基づく解析
 - STEP 4-3: 遺伝子オーソログに基づく解析



von Mering, C. et al. Science 2007, 31, 1126-1130

微生物によってゲノムサイズや有する遺伝子ファミリ、その遺伝子ファミリの コピー数や遺伝子長などは異なる。 このため、そのまま検索結果に表れた回数などから、相対存在度を 推定することはできない。



ほぼ全ての微生物が共通に1つだけ有する遺伝子ファミリをマーカーとして、 それらが検索結果に表れた回数から遺伝子長の補正を加えて 相対存在度の推定を行う



その他、分岐群特異的なマーカー遺伝子を用いる方法もある

分類群の相対存在度(1)





分類群の相対存在度(2)



頬粘膜

phylum_ratioの数値を利用して Excelなどで作図する



歯肉縁上の歯垢



アジェンダ



- GHOST-MP実行のながれ
- STEP 1: 実習の準備
 - STEP 1-1: SCLS計算機システムへのログイン
 - STEP 1-2: 必要なファイルのコピーと展開
- STEP 2: ジョブスケジューラの利用方法
- STEP 3: GHOST-MPによる配列相同性検索
 - STEP 3-1: GHOST-MPのコンパイル
 - STEP 3-2: データベースのインデックス作成(ghostmp_makedb)
 - STEP 3-3: 配列相同性検索(ghostmp_search)
- STEP 4: メタゲノムデータに対する相同性検索結果の解析
 - STEP 4-1: 検索結果の集計
 - STEP 4-2: 分類群に基づく解析
 - STEP 4-3: 遺伝子オーソログに基づく解析

遺伝子オーソログの相対存在度



配列相同検索の結果

summarize_search_result.py



ko_ratio ファイルの確認

[username@scls	buccal_mucosa]\$ sort -k 2 -gr result.buccal_mucosa.ko_ratio less
ko:K06147	0.003332
ko:K02029	0.003305
ko:K07052	0.003097
ko:K07024	0.002670
ko:K02003	0.002307
ko:K01992	0.002025
ko:K02028	0.002018
ko:K02006	0.001965
ko:K02030	0.001963
ko:K02004	0.001912
ko:K09687	0.001795
ko:K02529	0.001585
ko:K02913	0.001574
ko:K03574	0.001541
ko:K03402	0.001358
ko:K02078	0.001343
ko:K02008	0.001321
ko:K02015	0.001293
ko:K02793	0.001285
ko:K01462	0.001285

WebブラウザによるKEGG Orthologyの閲覧(1)



KEGG: Kyoto Encyclope.	× +				x
🗲 🔶 🛞 www.kegg.jp	▼ ♂ 【 Q、 検索	☆ 自	₩ 1	ø	Ξ
Experies	KEGG • K06147 Search Help » Japanese				•
KEGG Home Release notes Current statistics Plea from KEGG KEGG Database	KEGG: Kyoto Encyclopedia of Genes and Genomes KEGG is a database resource for understanding high-level functions and utilities of the biological system, such as the cell, the organism and the ecosystem, from molecular-level information, especially large-scale molecular databaset				н
KEGG overview Searching KEGG KEGG mapping	experimental technologies (See Release notes for new and updated features).				
最も割合の高かったk URL: http://www.keg	06147についてKEGGで調べる。 g.jp				
KEGG Software	KEGG PATHWAY KEGG pathway maps [Pathway list]				
KegTools	KEGG BRITE BRITE functional hierarchies [Brite list]				
KGML	KEGG MODULE KEGG modules [Module list Statistics]				
KEGG FTD	KEGG GENOME Genomes [KEGG organisms]				
Subscription	KEGG GENES Genes and proteins [Release history]				
Subscription	KEGG COMPOUND Small molecules [Compound classification]				
GenomeNet	KEGG REACTION Biochemical reactions [Reaction modules]				
Genomentee	KEGG DISEASE Human diseases [Cancer Infectious disease]				-
DBGET/LinkDB		http://w	/ww.kegg.jp	⊳/ ⇔ [1	/1] Тор

WebブラウザによるKEGG Orthologyの閲覧(2)



•) 🛞 www.g	enome.jp/dbget-bin/www_bget?ko:K06147	☆ 自 ♣	â	Ø	
K	ORTHOLOGY: K06147				
Entry	КО6147 КО	All links			
Name	ABCB-BAC	Ontology (7)			
Definitior	ATP-binding cassette, subfamily B, bacterial	KEGG BRITE (1)			
	ABC Transporters, Eukaryotic Type ABC8 (MDR/TAP) subfamily ABC8-BAC subgroup KO6147 ABC8-BAC; ATP-binding cassette, subfamily B, bacterial BRITE hierarchy	Gene (316454) KEGG GENES (11005) KEGG MGENES (304744) EGENES (207) OC (498) Protein sequence (9841)			
Other DBs	COG: COG1132 COG2274 COG5265 TC: 3.A.1.21 3.A.1.106 3.A.1.109	UniProt (9790) SWISS-PROT (51) All databases (326302)			
Genes	ISC: IscW_ISCW000765 EOJ: EC026_2861 EC026_2862 ECC: c2421 c2422 ECP: ECP_1939 ECP_1940 ECI: UTI89_C2180 UTI89_C2181(ybtP) ECV: APECO1_1055 APECO1_1056(ybtP) ECQ: ECED1_2246(irp) ECED1_2247(irp) ECK: EC55989_2204(irp) EC55989_2205(irp) ECT: ECIA139_1077(irp) ECIAI39_1078(irp) EOC: CE10_2258(ybtQ) CE10_2259(ybtP) >> show all	Download RDF			

http://www.genome.jp/



ここではiPATH2を用いて、検索結果に表れた 遺伝子オーソログが関連しているパスウェイを可視化する。 多くの遺伝子に関する情報を直観的に捉えることができる。

http://pathways.embl.de/



Yamada, T., et al. Nucleic Acids Research 2011, 39 (suppl 2), W412-W415

iPATH2によるパスウェイの可視化



iR iPath: Interactive Path × +				-			, 0	x
♦ ♦ ③ pathways.embl.de	⊽ C ^{el} 🛿 ▾ Google	٩	☆	ê	ŧ	⋒	9	≡
interactive Pathways Explorer 2			HOME	DATA HEL	.P TUTOR	IAL 8 ABC	DUT & CONT	LCT
Welcome to the Interactive Pathways Explorer v2								
Interactive Pathways Explorer (IPath) is a veb-based tool for the visualization, analysis and customization of the various pathways maps. Current version provides three different global overview maps:								
Select the devined version by driving the map receive below: Carbon Gration Gilyoxylate and Gilyoxylate and Gilyoxylate and metabolism Marking Gilyoxylate and Gilyoxylate and Marking Gilyoxylate and Marking	s based on KEGG data.							
Text zoom: 50%			http:/	//path	Ways.(embl.d	e/ ⇔ [1	/1] Тор

http://pathways.embl.de/にアクセスして、ここをクリック

iPATH2のインターフェース









ノードやエッジをクリックすると、化合物や反応に関する 情報を表示することができる

iPATH2の表示のカスタマイズ





「Customize」ボタンでパネルを表示してElement selectionを指定すると対応するパスウェイを強調表示できる。KEGG KOの例

iPATH2を用いたサンプル間の比較





頬粘膜と歯肉縁上の歯垢の結果をiPATH2を用いて比較する

例:割合が0.05%以上のKOの関わる代謝パスウェイをハイライトした

Path2: Metabolic overvi... x

Pathways.embl.de//Path2.cgi#
Customize
Search
Export
Metabolic pathways
Regulatory pathways
Biosynthesis of secondary metabolites
File of the secondary metabolite
File of the secondary metabolite
Metabolic pathways
Regulatory pathways
Regulatory pathways
Biosynthesis of secondary metabolites
File of the secondary metabolite
<p

頬粘膜

歯肉縁上の歯垢



本日の実習内容











監修・製作:東京工業大学 秋山研究室 GHOST-MPチーム 秋山 泰、石田 貴士、角田 将典、鈴木 脩司

資料製作 : (株)情報数理バイオ

ghost-mp@bi.cs.titech.ac.jp



2015年3月

独立行政法人理化学研究所 HPCI計算生命科学推進プログラム