Introduction to web based annotation system for Structural-Biological Whole Cell Project *Thermus thermophilus* HB8

Thermus thermophilus HB8用に開発された Web baseのアノテーションシステムの利用方法

Akinobu FUKUZAKI ¹,² 、 Fumikazu KONISHI ¹,² 、 Takeshi NAGASHIMA ¹ 、 Kaori IDE ¹,² 、 Mariko HATAKEYAMA ¹,² 、 Seiki KURAMITSU ²,³ 、 Akihiko KONAGAYA ¹,²

福崎 昭伸 1,2、小西 史一 1,2、長嶋 剛志 1、井手 香 1,2、 畠山 眞里子 1,2、倉光 成紀 2,3、小長谷 明彦 1,2

(¹Bioinformatics G., RIKEN GSC ²RIKEN Harima Institute ³Grad. Sch. of Sci., Osaka Univ.) (¹理研 GSC ゲノム情報科学研究G ²理研 播磨 ³大阪大学 理学研究科 生物化学)

e-mail: akki@gsc.riken.jp

高度好熱菌丸ごと一匹プロジェクトでは構造機能解析をはじめとして系統的機能解析のためのトランス クリプトーム解析やプロテオーム解析、メタボローム解析、システム生物学的研究を多数の研究者の協力 のもと進めている。その中でプロジェクト参加者の間で実験に必要なゲノム配列や遺伝子情報、タンパク 質の情報を共有する場を提供するとともにそれらの情報を更新できるようなシステムの必要性がもとめら れてきた。我々はWeb技術とデータベース技術を用いてプロジェクトに参加している研究者達が便利に活 用できるようなアノテーションシステムを指向して開発を続けてきた。今回は連携研究会の場をおかりし て我々が開発したシステムの使用方法の紹介を行う。

アノテーションシステムでできること

本システムでは現在のところ以下のような機能が利用できる。これらの機能は今後も拡張や追加がなさ れるのでより新しい説明はアノテーションシステムのWebサイトにあるオンラインドキュメントを参照し ていただきたい。





アノテーションシステムはユーザー認証システムによって保護されたWebサイトに設置されている。 よってシステムへ接続するにはあらかじめ申請手続きをしてユーザーアカウントを取得する必要がある。 本システムへのアクセスにはWeb Browserを利用するが、IEやMozilla等の一般的なWeb Browserで利用 が可能である。なるべく最新版を利用することをおすすめする。

https://access.obigrid.org/thermusへアクセスすると以下のようなユーザー認証の画面が表示される。 ここで取得したユーザー名とパスワードを入力しLOGINボタンをクリックする。

Open Bioinformates Grid	
A	ccess.obigrid.org ログイン
	ユーザー名 パスワード
	Locan
溆	f規ユーザ登録は <u>こちら</u> からどうぞ。
Secure Site cirk to write	Access.obigrid.org is fully contributed by <u>RIKEN-GSC</u> Copyright(C) 2002-2004 OBIGrid, All Rights Reserved. u have any question, please contact: info@obigrid.org

正常にログインできると以下のようにポータルサイトのFrontPageが表示される。ポータルサイトでは アノテーションシステムへのリンクだけでなくオンラインドキュメントやメンテナンス情報、さらには ユーザーが自由に使えるMy pageなどが提供される。ここではアノテーションシステムの説明を行う。ア ノテーションのメニューページへいくにはContentsのところのAnnotation Systemをクリックする。





アノテーションシステムはいくつかのメニューを持っている。現在のバージョンでは情報を閲覧する viewが1つとアノテーションを行うworkflowが4つある。これらのメニューをクリックすることでアノ テーションシステムの各機能へアクセスすることができる。以下に各機能の説明をする。

A	Annotation System https://access.obigrid.org/thermus/?Annotation%20System					
[New	Edit Freeze Diff Upload] [My page Front page					
4 ents /com	 Ver04 ORF/Genome view Function name annotation workflow GO Term annotation workflow PDB annotation workflow TTid annotation workflow 					
	FrontPage					

Ver04 ORF/Genome view

情報の閲覧を目的としている場合はこの機能を使用する。ORFをベースとした各種検索とアノテーション情報や配列情報、ゲノムのLocusの情報が提供される。

Function name annotation workflow

ORFの機能名を更新する時に使用する。

GO Term annotation workflow

GO Termの対応付けを更新する時に使用する。

PDB annotation workflow

PDBの対応付けを更新する時に使用する。

TTid annotation workflow

TTナンバーの対応付けを更新する時に使用する。ほとんどのTTナンバーORFは最新のゲノム配列に対 して位置を特定できており対応するORFも判明しているが、一部の例外がある。もし調べたいTTナンバー を検索しても見つからない場合はこちらで情報を更新する必要がある。



ここからは実際にシステムを使った例を示す。

TTナンバーでORFを検索する 🛛 キーワードで検索する 🖌 PDBで検索する

Ver04 ORF/Genome viewを使って各種条件で検索を行い、 該当するORFを呼び出すことができる。まずAnnotation SystemのメニューでVer04 ORF/Genome viewをクリックす る。すると右のような閲覧画面が表示される。

赤枠で囲まれたところが各種表示画面の切り替えタブで、ク リックすることで切り替えることができる。検索をしたい場合 はRegionSearchをクリックする。するとタブの下のツール表 示部分が検索用の画面に切り替わる。

青枠で囲まれたところがRegionSearchツールの中でさらに 機能を選ぶタブである。ここではまずSimpleを選ぶ。Simple ではTTナンバーやPDBのID、機能名のキーワードなどを入力 して検索することができる。

例えばTT0051とKeywordに入力してSearchをクリックする と2件の結果が表示される。TTHB84orF1388300という RegionはゲノムのForward strand側にあるORFで、もう一方 のTTHB84tt_5005100はTTナンバーORFをこのシステムの データベースに登録する時に付けられたIDである。このよう にTTナンバーORFとゲノムのORFの両方が見つかった場合は そのTTナンバーの対応付けがきちんとされていることにな る。検索結果のTTHB84orF1388300をクリックするとその ORFの情報を表示するモードに切り替わり、例えばHomolog のnr/ntタブをクリックすることでBLASTによりどのような配 列がヒットしたのかを確認することができる。ここでgi:のと ころをクリックするとNCBIのサイトへジャンプしてそのエン トリーの詳細を見ることができるし、ヒット位置を示すバーを クリックするとBLAST2Sequenceによる詳細なアライメント を見ることができる。

さらにPDBのIDで検索する時は11W7のようにPDB IDを keywordに入力することで該当するORFを見つけることがで きる。この場合は3件ヒットした。

また、DNA repairのようにスペースを間にいれて複数の キーワードを入力することでそれらの全てのキーワードが機能 名に含まれるORFを見つけることができる。この例ではDNA repair protein関係のORFだけでなくDNA mismatch repair proteinも見つかっている。



YourlD: akki Top / Annotation menu / Local BLAST



YouriD: akki Top / Annotation menu / Local BLAST TTHB8 Region Information Viewer TTHB840rF138830









配列のホモロジーで検索する

配列のホモロジーでORFを探索することもできる。例えば e.coliのLactose Operon repressor (P03023)をNCBIのサイト からとってきて検索してみる。RegionSearchツールのBLAST タブをクリックすると配列を入力する画面になる。ここに P03023のペプチド配列をペーストし、Sequence Typeを Peptideにセットしてから、searchをクリックする。

するとこの配列と似た配列を持つORFがいくつか表示され る。このように配列のホモロジーでORFを探す機能もあるので 機能未知の配列であったり、まだアノテーションがきちんとさ れていない配列であってもゲノム上のどのORFに近いものかを 見つけることができる。

Sequence:

Sequence Type: ONucleatide Pentide
V OBES on Genome VITI OBES
(search)
view results by TSV

PDB strand pos_start pos_end dna_len ref_id (Option)						
reg_id	ttid	ann_name	percent_identity	evalue	score	
HB84tt_5106000	TT1060	transcriptional regulator, LacI family	28.79	2.9e- 24	106.3	
HB84orF2015000	TT1060	facl-family transcriptional regulatory protein	28.79	3.8e- 24	105.9	
HB8411 5104000	TT1040	probable transcriptional regulator	28.8	7.5e- 20	91.66	
HB8411_5101600	TT1016	transcriptional regulator	25.07	2.3e- 13	70.09	
HB84orR1094400	TT1016	transcriptional regulator	25.07	2.3e- 13	70.09	
HB84tt 5090700	TT0907	S-adenosylmethione:tRNA ribosyltranferase-isomerase	50.0	8.6	25.02	
HB84tt 5021500	TT0215	alpha-ribazole-5'-phosphate phosphatase (CobC)	29.27	8.6	25.02	
HB84orF2009200	TT0215	alpha-ribazole-5'-phosphate phosphatase	29.27	8.6	25.02	

TTナンバーの対応付けの更新 🛛 機能名の更新 🖌 PDBの対応付けの更新

アノテーションデータはまだ多くの研究者らによる確認を必要としており、我々のシステムではアノテーションデータを更新する機能も提供している。ここではPDB IDの対応付けのアノテーションを例にその手順を示す。

Annotation SystemメニューでPDB annotation workflowを クリックすると右のような画面が表示される。Workflowは多 数の配列に対するアノテーションを効率的に進めるためのフ レームワークで、画面左上のWorkFlowのところにあるよう に、3つのフェーズからなるユーザーインターフェースで構成 されている。

まずはターゲットとなる配列群を探し出すRegion Search フェーズで、閲覧の時に使用したRegionSearchと同様の画面 を用いて検索する。事前に処理してあるホモロジーの結果から アノテーションするのでここではRevHomologのタブをクリッ クして対象の配列を探す。

例としてBiosynthesis Enzymeのキーワードが入っている PDBエントリーにヒットした配列を探してみる。Keywordに Biosynthesis Enzymeと入力し、PDB HitのHitを選択、 percent identityをLE (Lower Equal)を選択して60.0と入力 してSearchをクリックする。

結果として1UC8と1UC9のエントリーが見つかる。これら をRegionListに登録するためにadd all regionsをクリックす る。ターゲットとなる配列のIDはいずれも同じものだったた め、重複が自動的に判別されてひとつのIDがRegionListにいれ られる。

これでアノテーションの準備ができた。検索ツールのすぐ下 にあるnextボタンをクリックしてPDB annotationのフェーズ へうつる。





4 entries:					
PDB strand pos_start pos_end dna_len ref.					
reg_id	hit_id	percent_identity	score	evalue	
add TTHB84orR1312300	1UC8	95.0	498.8	0.0	
add)TTHB84orR1312300	1UC8	95.0	498.8	0.0	
add TTHB84orR1312300	1UC9	94.29	495.4	0.0	
add TTHB84orR1312300	1UC9	94.29	495.4	0.0	
add all regions					



PDB annotationフェーズは右の画面のような構成になって いる。左側は先ほどのRegion Searchフェーズと同じ構成で WorkFlowのところはPDB annotationフェーズのところに色 がついて表示されている。右側はPDB IDの対応を入力するア ノテーションインターフェースで、その下は閲覧の時と同じよ うな配列情報を表示するツールになっている。右の画面の例で はHomologsタブをクリックしPDBタブをクリックしてPDBエ ントリーとのホモロジー情報を表示している。

Homologsツール中ではターゲット配列とホモロジーが見ら れたエントリーのより細かい情報が見られるようになっている がデフォルトではいくつかの情報が表示されていない。例えば percent identityとhit len(アライメントがとられた範囲の配 列の長さ)を同時に見たい場合はホモロジー情報の上の各項目 のチェックボックスにチェックをしてOptionをクリックする ことでこれらの情報が表示される。この例ではORFの長さが 843bpなのに対してhit lenが840なので1コドンだけ足りない らしいことがわかるし、ホモロジーのバーの右端に少し黒い帯 が見えていることから終端のコドンが外されているようにも見 える。

アライメントを詳細に見たい場合はこのバーをクリックする ことで詳細なアライメントを見ることができる。アライメント はBLAST2Sequenceを使ってリアルタイムに計算されるので 必要であれば各種パラメーターを調整して再計算させることも できる。

ほかにもホモロジー情報のところのgi:のところはNCBIの Entrezにリンクしてあってこのエントリーの詳細を見ることが できるし、pdb:のところはPDBのデータベースにリンクして あってより詳細な情報をおいかけていくことができる。

このようにしてホモロジー情報を確認して当該のPDB IDが このORFに付けられるべきであると判定できたらアノテーショ ンインターフェースの空欄にPDBやChain名など必要事項を入 力してaddをクリックすることでPDBアノテーションを追加す ることができる。またアノテーションインターフェースの上に すでに付けられたアノテーション情報が表示されているのでこ れを見て確認することもできるし、PDBの横にあるDelボタン をクリックすることでこのアノテーションをORFから外すこと ができる。アノテーションがすんだらnextをクリックすること で次の配列へ移動するが、この例では配列は一つだったので Reportフェーズへうつり、今回のアノテーション結果をまと めて確認できる。

アノテーション情報は即座にデータベースへ反映されるため、この後すぐに閲覧の画面等でRegion Searchを使って先ほどのPDB IDを検索するとこの配列のIDが見つかる。







uery:	1	MLAILYDRIRPDERMLFERAEALGLPYKKVYVPALPMVLGERPKELEGVTVALERCVSQS	180
bjct:	1	MLAILYDRIRPDERMLFERAEALGLFYKKVYVPALPMVLGERPKELEGVIVALERCVSQS	60
uery:	181	RGLAARYLTALGIPVVNRPEVIEACGDKWATSVALAKAGLPQPKTALATDREEALRLME	360
bjct:	61	RGLAAARYLTALGIPVVNRPEVIEACGDKWATSVALAKAGLPQPKTALATDREEALRLME RGLAAARYLTALGIPVVNRPEVIEACGDKWATSVALAKAGLPQPKTALATDREEALRLME	120
uery:	361	AFGYPVVLKPVIGSWGRLLAKVTDRAAAEALLEHKEVLGGFQHQLFYIQEYVEKPGRDIR	540
histe	121	AFGYPVVLKPVIGSWGRLLA KEVLGGFQHQLFYIQEYVEKPGRDIR	190
uery:	241	VEVYORALAALIAALIAADAAWIINIAACOQAENCELTEEVARDSVKAAEAVGGGVVAVDLFE	720

gi:37928151 Chain A, Crystal Structure Of A Lysine Biosynthesis Enzyme, Lysx, From Thermus Thermophilus Hb8.	PID 95.00
gi:37928152 Chain B, Crystal Structure Of A Lysine Biosynthesis Enzyme, Lysx, From Thermus Thermophilus Hb8.	PID 95.00%
PDB id: 1UC9 Chain: A	
GI: 37928155	
Name: Chain A, Crystal Structure Of A Lysine Biosynthesis E	
Comment: PID 94.29%	
add	





その他のツール

本システムでは他にも情報の閲覧やアノテーションを助ける ツールが用意されている。

Annotation Historyではアノテーションが付けられていった 経緯を見ることができるため、判定が難しいアノテーションで はどのような意見が過去に出てアノテーションされたかを見る ことができる。

FastaではNucleotide(DNA)とPeptide(AA)の配列をFasta形 式で取得することができる。BLAST等のツールへここでコ ピーした配列をペーストして使用することができる。Peptide の配列はORFが設定されている読み方向にたいして3段階のフ レームシフト全部を用意してある。

RegionViewではORFの配列の様子をみることができる。 DNA配列とAA配列の両方がゲノム上のマップで表示されてお り、他の領域とどのように重なっているのかなど詳細にみるこ とができる。また5'UTR、3'UTRのボタンを用意してあるので ORFの両末端付近の配列を素早く確認することもできる。さら に表示位置を移動したり、表示範囲を広げたりすることができ るため、ORFのところだけでなくさらに上流や下流をずっと 追っていくことができるようになっている。

Genome BrowseではView current RegIDをクリックするこ とでそのORFの近辺のゲノム情報を見ることができる。データ ベースに登録されているORFやORFの間のIntergenic領域がエ ントリーとして表示されそれらの配置がバーで表示されるた め、各配列の大きさや位置関係が直感的に把握できるように なっている。また機能名をリストで見やすく表示してあるので ターゲットの遺伝子の前後にどのような機能のものがあるかを 簡単に確認できるようになっている。このツールにも参照場所 をスクロールさせるボタンやズームイン、ズームアウトがある ので自由に表示領域を変更することができる。









まとめ、今後

遠隔地から閲覧できるアノテーションデータベースという課題を満たすため我々はWeb技術を用いて アノテーションシステムを構築した。今回は基本的なツールと基礎的なデータ構造やシステムアーキテ クチャの検討を行いながら実装をしていくことが目標であったが最低限のものは用意できたように思わ れる。しかし、"丸ごと一匹"を実現するためにはこれから先より多くの知見を効率的に蓄積し、利用で きるよう環境整備をすすめていかなくてはならない。特に代謝系の情報やMS、GeneChipなどの発現 データなど実験の結果などを集約していけるように速急に対応していく予定である。