

2008年1月8日

# ライフサイエンス統合データベースプロジェクト 共通基盤技術開発チーム進捗状況報告 (CBRC)

CBRC: 浅井、野口、諏訪、  
金、ホートン、広川、福井、藤渕、堀本

# 担当業務：ライフサイエンス統合データベース開発運用 (共通基盤技術開発)

## ワークフロー技術を用いた統合DB環境の構築

利用者が得たい情報(知識)を、パソコンなどの端末から要求すると、必要なデータ、解析手法などを、国内、海外から自動的に選び、データベースと解析ツールのワークフローを作成し、最適な計算資源を使って解析を行う統合DB環境を構築する。

# 統合データベースのイメージ



バイオインフォマティクス  
統合環境ウェブページ

**Integrated Environment for Bioinformatics Knowledge Discovery**

- Structure based drug design
- Target gene screening
- Side effect prediction

DB検索要求  
データ解析要求  
結果取得

DB検索要求  
結果取得

独自DB、ツールと  
大規模計算を組み合わせた  
解析要求・結果取得



海外サイト  
ソフト・DB

国内サイトのDB



例 この薬の  
副作用を予測したい

化合物  
+  
本来の標的遺伝子

最新のIT技術を利用し、  
必要なデータ、解析手法  
を国内、海外から自動的  
に選びパイプラインを構築、  
最適な計算資源を使って  
解析を実行して返す

相互作用・結合DB検索  
構造予測ベース解析  
ターゲット遺伝子予測  
細胞内ネットワーク解析



目的外標的遺伝子リスト  
分子設計指針  
副作用予測結果

データベース群

統合データベースセンター

ツール群

Mirror Site: DDBJ, PIR, PDB, SWISS-PROT, KEGG, etc.

計算資源



# 統合DBにおけるワークフローの開発方針

## • グリッドを利用

理由1: 目標とする統合DBを構築するために必要な技術が、ほぼグリッド環境で実現できている。(例: GEOGrid)

- 長時間JOBの対応
- ソフトウェアやデータベースの利用制限(アカデミック利用のみ等)
- その他

理由2: ワークフロー実現のために必要で、かつ現状のグリッド技術にない技術の開発に関しては、産総研グリッド研究センターの協力が得られる。

## • C B R C 内でプロトタイプを構築し、外部と連携

# 今年度の計画

- CBRC内の計算機環境で、固定のデータベースと解析ツールを用いたワークフローを作成し、統合DB環境のプロトタイプを構築する。

- **グリッド環境構築**

- **ワークフローの議論**

- **ソフトウェア・データベース間のインターフェイス**

共通基盤技術開発チーム内の打ち合わせで、SOAPに決める

**CBRCにおいて公開しているWebサービスのSOAP**

**ワークフロー固定**

**ワークフロー可変 (taveruna, KNIME調査)**

- **グリッドWebアプリケーションの開発**

# システム構成



産総研外ソフトウェアおよびDB

標準的なWeb Serviceで統合

OGSAによる標準アーキテクチャ

グリッド環境の計算資源  
への計算依頼

データグリッドによる  
大規模データの仮想化  
OGSA-WebDB

産総研内計算資源

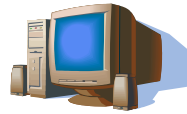
グリッド用計算サーバ、Magiクラスターなど

産総研内DB



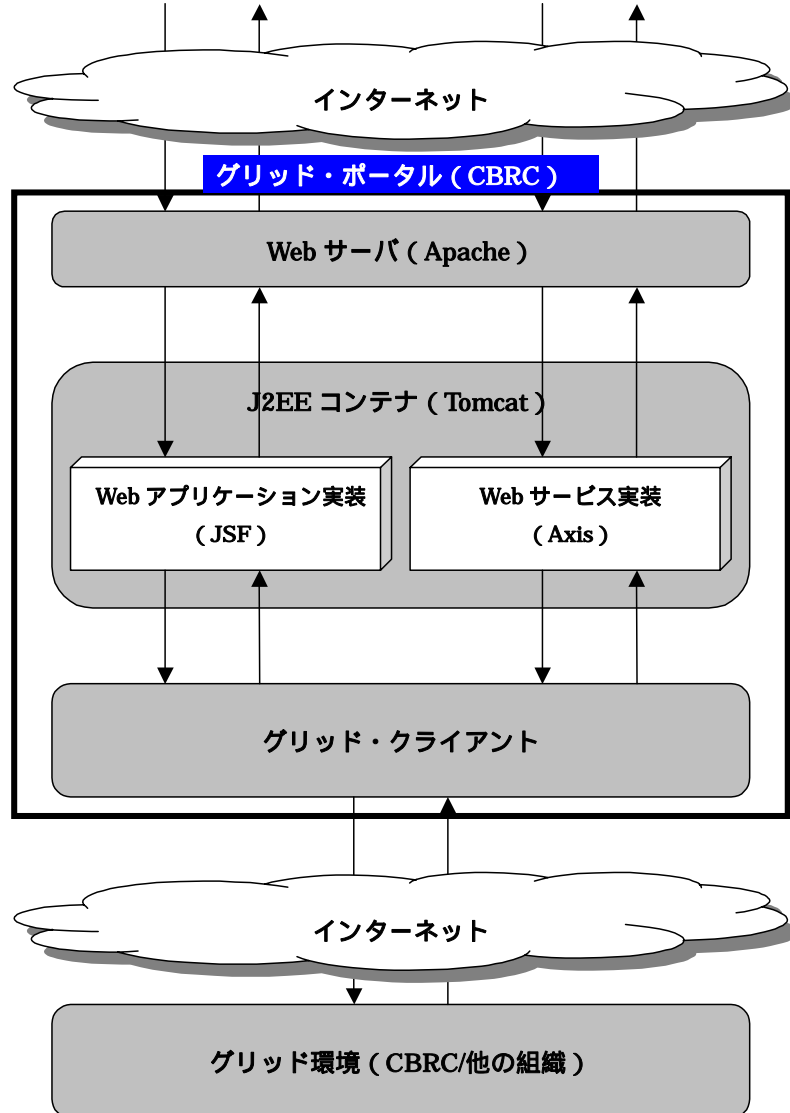


Web アプリケーション利用者  
(Web ブラウザ)



Web サービス利用者  
(SOAP クライアント)

# グリッド環境構築



ハードウェア:

- ・グリッド用Webサーバ
- ・グリッド管理サーバ
- ・計算サーバ(発注済)

ソフトウェア:

- ・OS(CentOS 5.0)
- ・Web関連
- ・グリッド関連

# CBRCで研究開発された公開ソフトウェア(例)

## Software

### ALN

2つの配列または配列グループ間をア  
<http://www.cbrc.jp/ALN/>

### ASIAN

階層クラスタリングとグラフィカル・ガ  
み合わせによるネットワーク推定ツ  
<http://eureka.cbrc.jp/asian/>

### CellMontage

マイクロレイデータ検索・解析システ  
<http://cellmontage.cbrc.jp/>

### CoCoozo

ペプチド・タンデム質量分析向け並列  
<http://www.cbrc.jp/cocoozo/>

### FORTE

プロフィール比較によるタンパク質立  
<http://www.cbrc.jp/forte/>

### GeneDecoder

真核生物ゲノムのDNA塩基配列から遺伝子を予測するソフトウェア  
<http://genedecoder.cbrc.jp/>



### GRIFFIN

GPCR-タンパク質結合選択性予測シ  
<http://griffin.cbrc.jp/>

### GUPPY

遺伝子配列におけるデータの意味を註釈  
グラム  
<http://www.cbrc.jp/GUPPY/>

### Murlet

構造RNAのマルチプルアライメントプロ  
<http://murlet.nornia.org/>

### PAPIA

並列タンパク質情報解析システム  
<http://mbs.cbrc.jp/papia/>

### POODLE

タンパク質ディスオーダー予測  
<http://mbs.cbrc.jp/poodle/>

### Scarna

類似RNA高速検索ツール  
<http://www.scarna.org/>



### SGCAL

糖鎖断片化解析ツール  
<http://sgcal.cbrc.jp/>



### SOKOS/CAN

多目的RNA配列解析ツール  
<http://www.cbrc.jp/sokos/>



### TMBETA-NET

膜貫通タンパクのベータストランド予測をするプログラム  
<http://psfs.cbrc.jp/tmbeta-net/>



### WoLF PSORT

タンパク質細胞内局在化予測ソフト  
<http://wolfsort.org>





# その他の主な公開ソフトウェア

- RNAmine <http://www.scarna.org/rnamine/> (SOAP化済)  
~ RNA配列からの2次構造モチーフ探索サービスを提供している。複数のRNA配列(構造アノテーション不要、アラインメント不要)を入力すると、そこから2次構造モチーフを抽出する。独自のアルゴリズムによって、極めて高速なモチーフ抽出を実現しているところが特徴。
- Murlet <http://murlet.ncrna.org/> (SOAP化済)  
~ RNA配列用高速多重配列アラインメントサービスを提供している。複数本のRNA配列(構造アノテーション不要)を入力すると2次構造を考慮した配列アラインメントの結果を出力する。独自の工夫によって、従来手法に比べて高速な処理を実現しているところが特徴。
- PHMMTS <http://www.scarna.org/phmmts/> (SOAP化済)  
~ RNA配列用構造アラインメントサービスを提供している。構造アノテーション付RNA配列(複数本可)と、構造未知のRNA配列1本を入力として、構造未知のRNA配列に対する構造アノテーションを推定する。RNA2次構造情報を木構造と隠れマルコフモデルの組合せで表現するという新しいアイデアを実証するための proof of concept 的サービス。
- MXSCARNA <http://mxscarna.ncrna.org/>  
~ RNA配列用超高速多重配列アラインメントサービスを提供している。複数本のRNA配列(構造アノテーション不要)を入力すると2次構造を考慮した配列アラインメントの結果を出力する。独自のアルゴリズムによって、従来手法に比べて極めて高速な処理を実現し、かつ精度も同程度を維持しているところが特徴。

# CBRCで研究開発された公開データベース(例)

## Databases

### ASTRA

選択的スライシング・選択的転写開始パターン分類データベース  
<http://alterna.cbrc.jp/>

### ConfC

タンパク質構造変化部位データベース  
<http://mbs.cbrc.jp/ConfC/>

### DB-SPIRE

タンパク質機能構造データベース  
<http://mbs.cbrc.jp/DB-SPIRE/>

### EzCatDB

酵素触媒機構データベース  
<http://mbs.cbrc.jp/EzCatDB/>

### fRNA Database

fRNAdbとUCSC GenomeBrowser for Functional RNAで構成される  
新規機能性RNA遺伝子発見支援用データベース  
<http://www.norna.org/>

### GENIUS II

全ゲノムタンパク質の立体構造帰属データベース  
<http://genius.cbrc.jp/>



### INOH

シグナル伝達パスウェイ代謝パスウェイデータベース  
<http://www.inoh.org/>



### PDB-REPRDB

配列、および構造類似性を考慮した代表タンパク質チェーンデータ  
セットを作成するシステム  
<http://mbs.cbrc.jp/pdbreprdb/>



### SEVENS

GPCR遺伝子の網羅的データベース  
<http://sevens.cbrc.jp/>



### TMBETA-GENOME

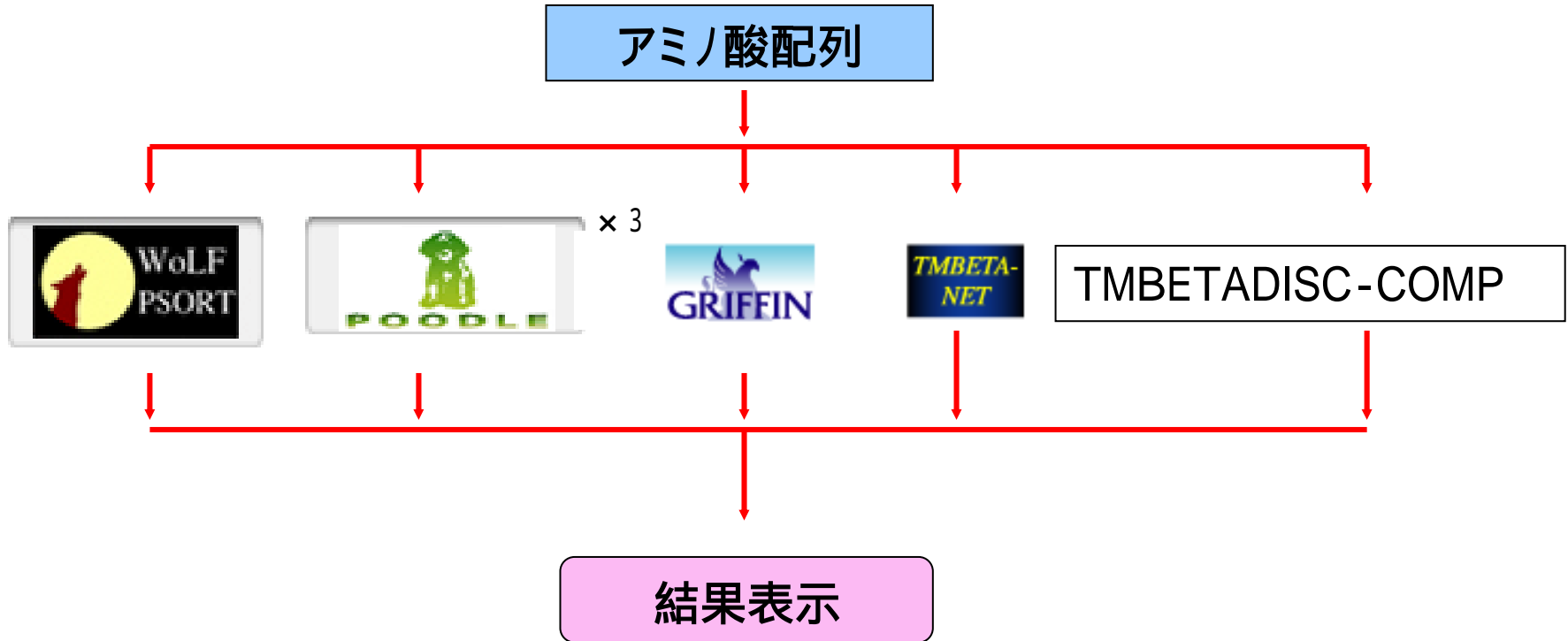
ゲノム中に含まれる $\beta$ -平行型膜タンパク質のデータベース  
<http://tmbeta-genome.cbrc.jp/annotation/>



# CBRC内公開サービスのSOAP化状況


	計
SOAP化されたDB	0
SOAP化されたソフトウェア	4
SOAP化中のDB	0
SOAP化中のソフトウェア	7
SOAP化未対応のDB	10
SOAP化未対応のソフトウェア	15

# プロトタイプ1 (配列解析)



注) DBCLSの予算でSOAP化

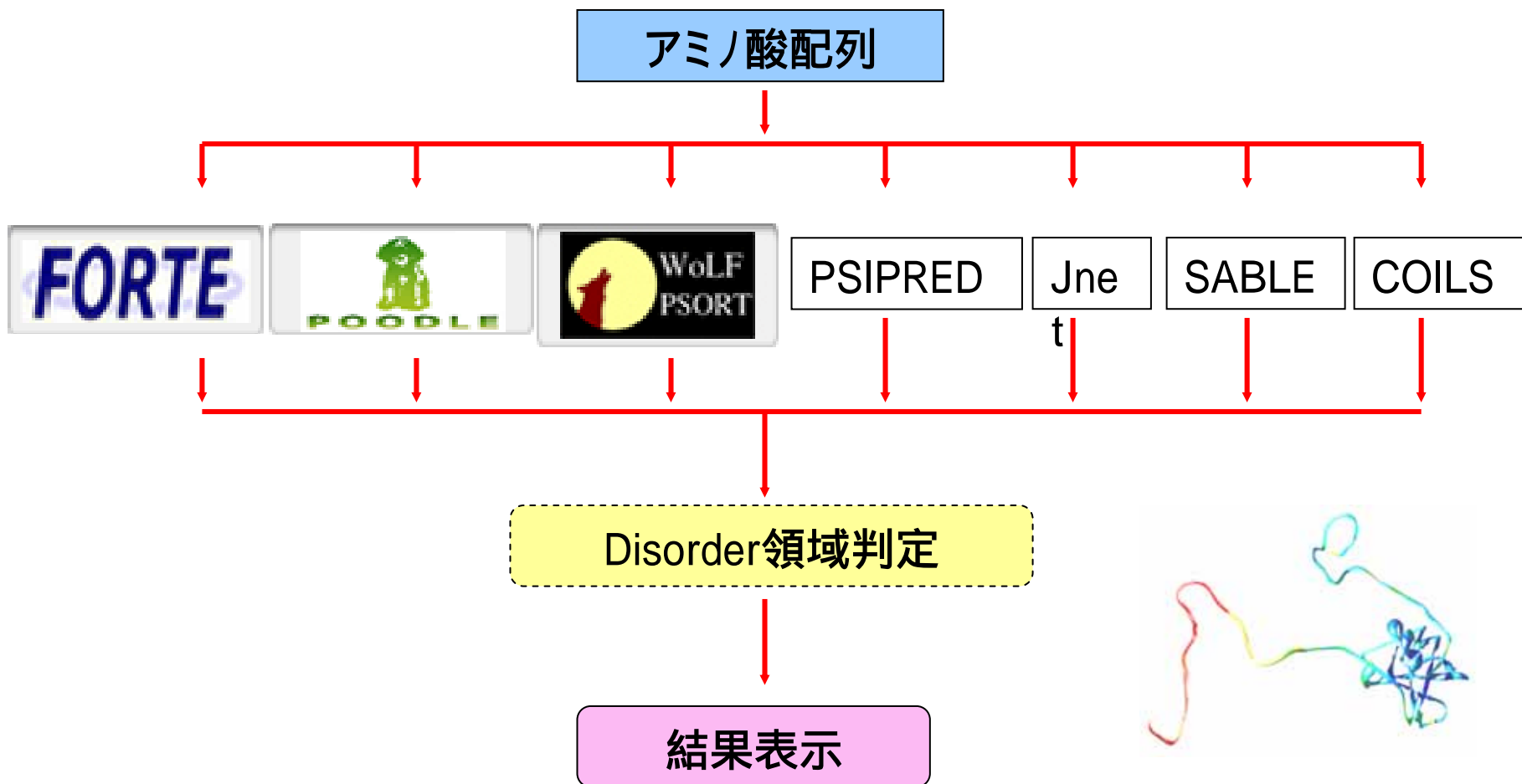
# プロトタイプ結果表示イメージ

		10	20	30	40																														
<b>Disorder prediction</b>																																			
Disorder領域予測	POODLE-S	<table border="1"> <tr> <td>MANLGCWMLV</td> <td>LFVATWSDLG</td> <td>LCKKRPKPGG</td> <td>WNTGGSRYPG</td> </tr> <tr> <td>1111122233</td> <td>2233344555</td> <td>6678898877</td> <td>7766798999</td> </tr> <tr> <td>MANLGCWMLV</td> <td>LFVATWSDLG</td> <td>LCKKRPKPGG</td> <td>WNTGGSRYPG</td> </tr> <tr> <td>1112232333</td> <td>4455666777</td> <td>7789899997</td> <td>7889999999</td> </tr> <tr> <td>MANLGCWMLV</td> <td>LFVATWSDLG</td> <td>LCKKRPKPGG</td> <td>WNTGGSRYPG</td> </tr> <tr> <td>1112232333</td> <td>4433455666</td> <td>8888899999</td> <td>9999999999</td> </tr> <tr> <td colspan="2">SCORE: 0.9476,</td> <td colspan="4">Disordered.</td> </tr> </table>				MANLGCWMLV	LFVATWSDLG	LCKKRPKPGG	WNTGGSRYPG	1111122233	2233344555	6678898877	7766798999	MANLGCWMLV	LFVATWSDLG	LCKKRPKPGG	WNTGGSRYPG	1112232333	4455666777	7789899997	7889999999	MANLGCWMLV	LFVATWSDLG	LCKKRPKPGG	WNTGGSRYPG	1112232333	4433455666	8888899999	9999999999	SCORE: 0.9476,		Disordered.			
	MANLGCWMLV					LFVATWSDLG	LCKKRPKPGG	WNTGGSRYPG																											
	1111122233					2233344555	6678898877	7766798999																											
	MANLGCWMLV					LFVATWSDLG	LCKKRPKPGG	WNTGGSRYPG																											
	1112232333					4455666777	7789899997	7889999999																											
MANLGCWMLV	LFVATWSDLG	LCKKRPKPGG	WNTGGSRYPG																																
1112232333	4433455666	8888899999	9999999999																																
SCORE: 0.9476,		Disordered.																																	
POODLE-L																																			
POODLE-T																																			
POODLE-W																																			
<b>Secondary structure prediction</b>																																			
二次構造予測	PSIPRED	<table border="1"> <tr> <td>MANLGCWMLV</td> <td>LFVATWSDLG</td> <td>LCKKRPKPGG</td> <td>WNTGGSRYPG</td> </tr> <tr> <td>9611188999</td> <td>5576778761</td> <td>5567647384</td> <td>9878989999</td> </tr> <tr> <td>MANLGCWMLV</td> <td>LFVATWSDLG</td> <td>LCKKRPKPGG</td> <td>WNTGGSRYPG</td> </tr> <tr> <td>9611188999</td> <td>5576778876</td> <td>1888888888</td> <td>8888888888</td> </tr> <tr> <td>MANLGCWMLV</td> <td>LFVATWSDLG</td> <td>LCKKRPKPGG</td> <td>WNTGGSRYPG</td> </tr> <tr> <td>1112232333</td> <td>4433488888</td> <td>8888888888</td> <td>8888888888</td> </tr> </table>				MANLGCWMLV	LFVATWSDLG	LCKKRPKPGG	WNTGGSRYPG	9611188999	5576778761	5567647384	9878989999	MANLGCWMLV	LFVATWSDLG	LCKKRPKPGG	WNTGGSRYPG	9611188999	5576778876	1888888888	8888888888	MANLGCWMLV	LFVATWSDLG	LCKKRPKPGG	WNTGGSRYPG	1112232333	4433488888	8888888888	8888888888						
	MANLGCWMLV					LFVATWSDLG	LCKKRPKPGG	WNTGGSRYPG																											
	9611188999					5576778761	5567647384	9878989999																											
	MANLGCWMLV					LFVATWSDLG	LCKKRPKPGG	WNTGGSRYPG																											
9611188999	5576778876	1888888888	8888888888																																
MANLGCWMLV	LFVATWSDLG	LCKKRPKPGG	WNTGGSRYPG																																
1112232333	4433488888	8888888888	8888888888																																
Jnet																																			
SABLE																																			
膜貫通領域予測 																																			
細胞内局在予測	<i>k used for kNN is: 27</i>																																		
	Q8TC27 ADA32_HUMAN <a href="#">details</a> E.R: 7.0, plas: 6.0, pero: 5.0, golg: 5.0, extr: 2.0																																		

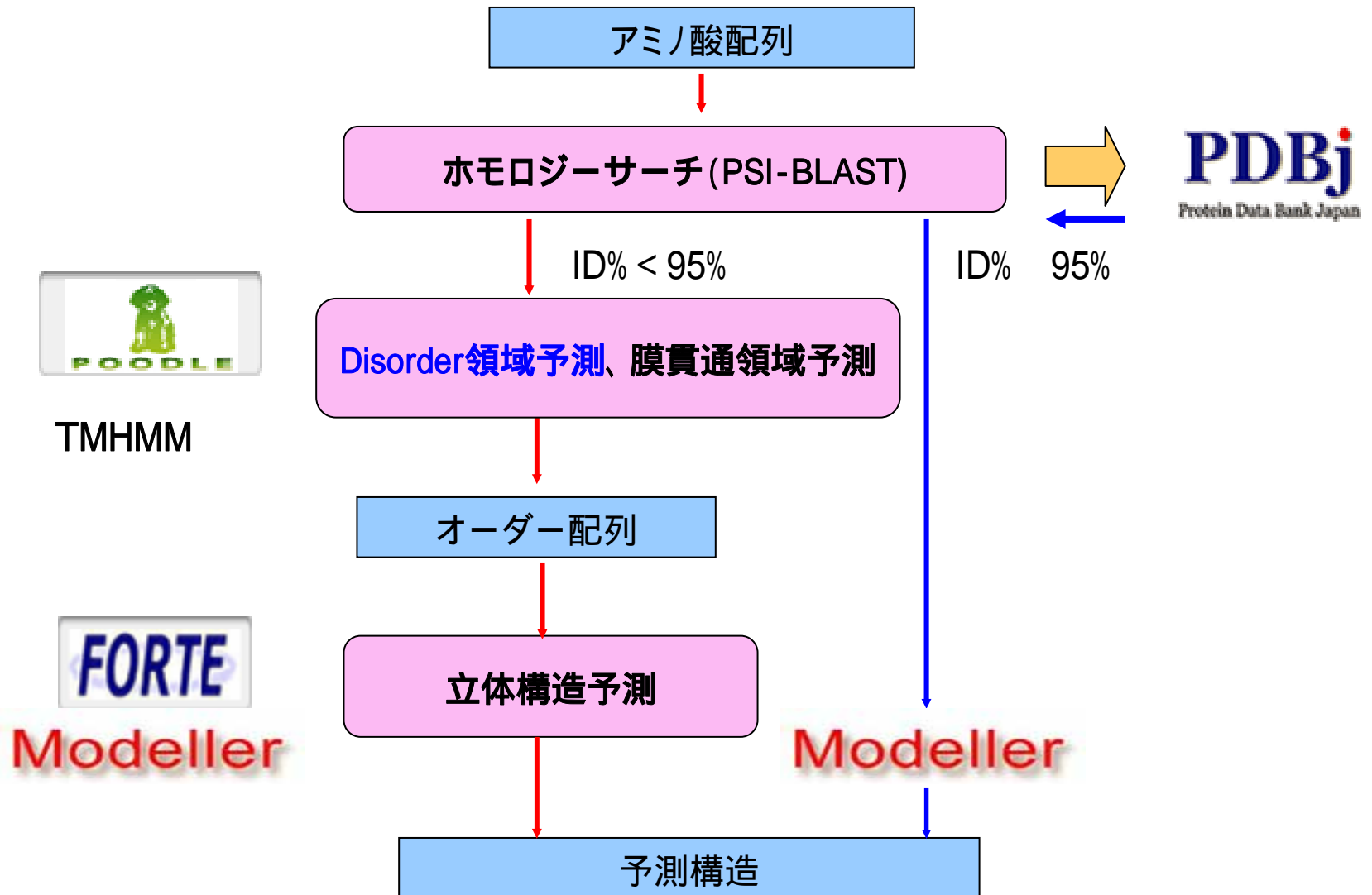
CBRCの公開Webサーバーの結果をこのイメージで貼り付けていく

# プロトタイプ2

( CASP7で用いたDisorder領域予測手法 )  
~ 総合評価で2位、長いDisorder領域で1位 ~

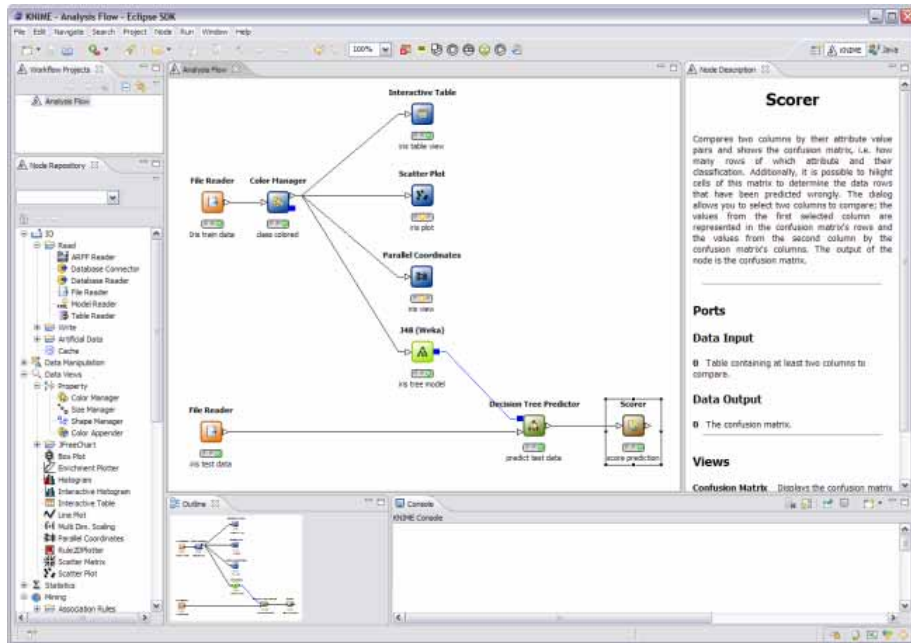


# タンパク質立体構造予測ワークフロー

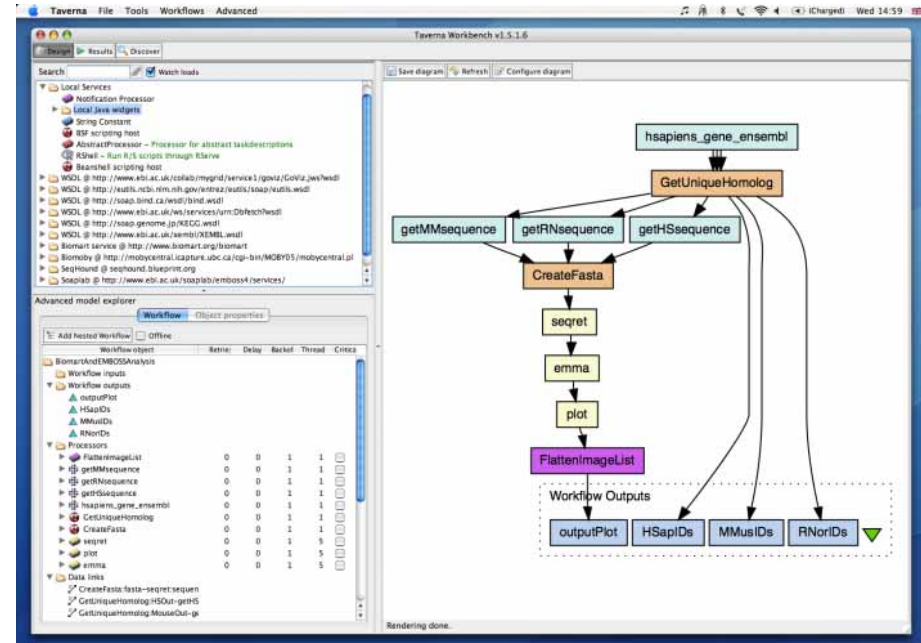


注) JST BIRD「タンパク質の構造・機能予測法の開発とヒトゲノム配列への適用」(代表: 東工大 太田元規准教授) において研究開発中

# ワークフロー技術調査



KNIME



taveruba



# 今後の予定

## • DBCLSとの連携

- 1) CBRC のサービスが DBCLS のサーバから利用可能
- 2) CBRC のサービスにDBCLSのサーバのサービスを組み込む
- 3) 上記以外の連携は、定期的な打ち合わせを行い議論

## • CBRC内のグリッド環境

- 1) DBCLSとグリッド環境の連携
- 2) グリッド用計算サーバ整備

## • ワークフローのサービス

- 1) 他のサービスの検討・設計
- 2) ワークフローを構成するソフトウェアのモジュール化  
(taveruna、KNIMEのイメージ)

# 開発計画

	2007	2008	2009	2010
<b>グリッド環境構築</b>	基本環境 —————→	必要に応じて機能拡張 —————→		
<b>プロトタイプ開発</b>	—————→			
<b>実用サービス検討・開発</b>		—————→	—————→	—————→
<b>DBCLSとかずさとの連携の検討・開発</b>		—————→	—————→	—————→