



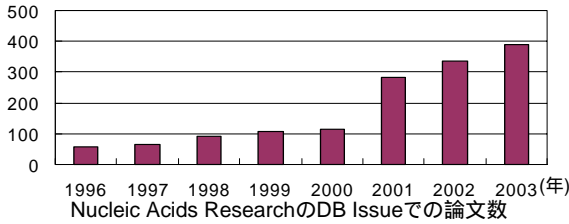
DataGrid: グリッド技術によるバイオ データベースの連携

大阪大学 情報科学研究科
松田秀雄

生命科学がグリッドを必要とする理由

- 生命科学はゲノムの登場によってdata driven/intensive scienceに変わりつつある
- 原子、分子、細胞、組織、さまざまなレベルでのシミュレーションが可能であるが、**大規模な計算パワー**を必要とする
- 生命科学は歴史もあり、非常に多岐にわたっているため、そもそも**科学者の協調**が必要である
- 生命科学にまつわる大量の情報が世界中に分散するデータベースで維持されており、それらの**連携**が必要

● ゲノムプロジェクトの成果によるデータベースの数および量の急激な増大



Domain	No. of DBs
DNA	87
RNA	29
Protein	94
Genomic	58
Mapping	29
Protein structure	18
Literature	43
Miscellaneous	153

GenBankのデータ量の増大
<http://www.ncbi.nlm.nih.gov/Genbank/genbankstats.html>

Current total : 511 Bio DB Catalog (DBCAT)への登録数 : <http://www.infobiogen.fr/services/dbcatal/>

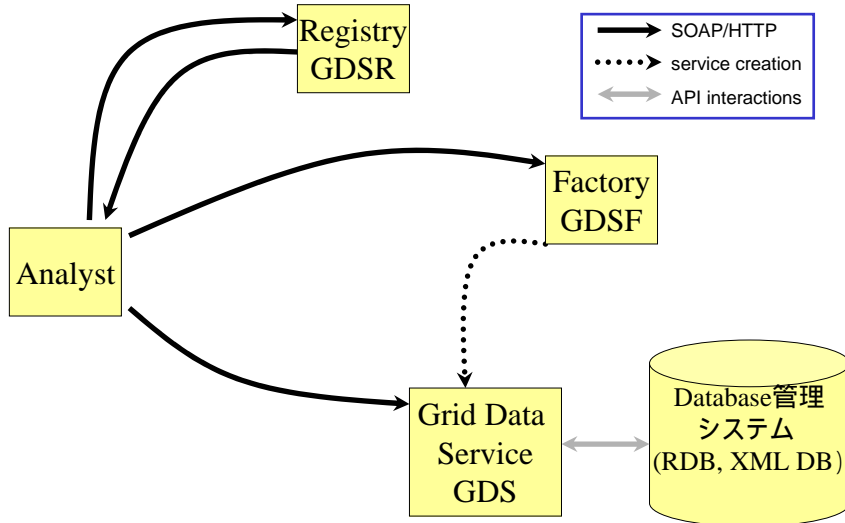
多数のデータベースの利用法

- ハイパーリンク結合 (例: 京大DBGET / LinkDB)
 - 既に多数のデータベースを連携する方法として実用的に運用されている。
 - ハイパーリンクに「意味」を持たせられない。
- 統合データベース (例: NCBIのEntrez)
 - 個々のデータベースを意識する必要がない。
 - 元のデータベースのスキーマの更新があると、統合のやり直しが必要。
- 異種データベース (例: スタンフォード大学のTSIMMIS)
 - データベースごとにラッパーを作成して共通形式へ動的に変換し、メディアエータで相互の検索を仲介 (スキーマ更新の影響は1つのラッパーのみ)
 - 生命科学特有のデータベース検索 (例えば、配列ホモロジー検索や構造類似検索) と組み合わせることが困難。
 - 商用データベースの利用でのユーザ認証への対応が困難。

グリッドによるデータベースの連携(データグリッド)

- データ検索の方法 (キーワード検索、全体検索など) ごとにグリッドサービスを提供し、サービス連携によりデータベース間の相互参照を行う。
- グリッドの認証機構を利用して、各データベースの認証を一元化。

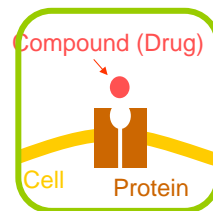
OGSA-DAI (Open Grid Service Architecture
Data Access and Integration) の利用



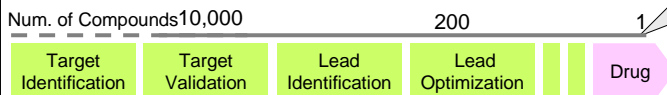
Drug Discovery Process

Application to the **drug discovery process**.

- ◆ Compounds (drugs) are activated by binding to proteins in a cell.
- ◆ **Drug Discovery Process** is to find chemical compounds that have good effects on their target proteins.
- ◆ The process needs **much time and much money**.

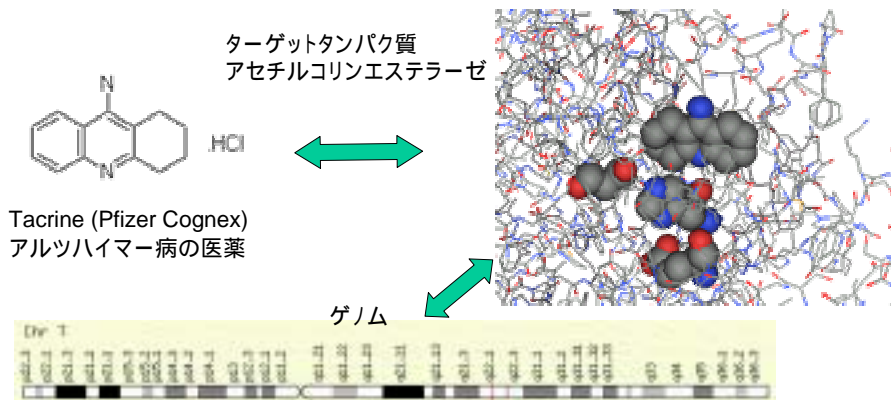


5~10 million \$
(10~15 years)



Protein-Compound Interaction Search is one of the most important technologies in drug discovery.

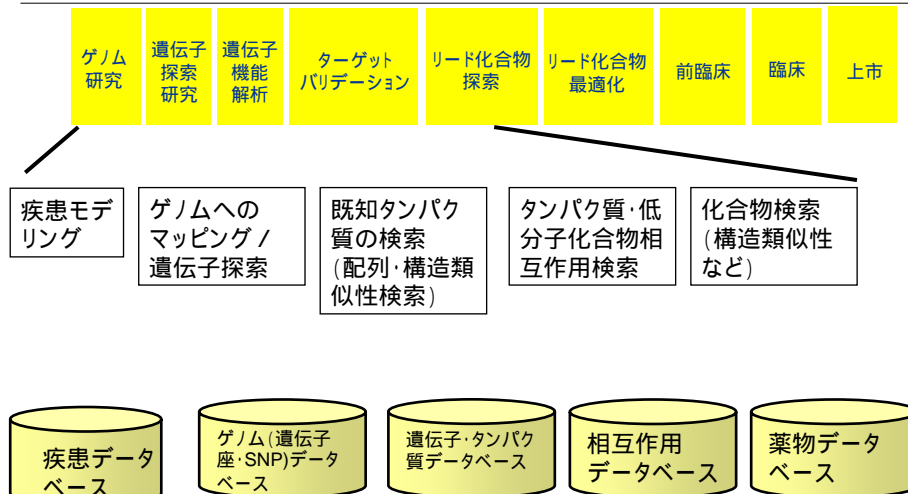
DataGrid can bridge **Biology** (proteins) and **Pharmaceutics** (compounds) !

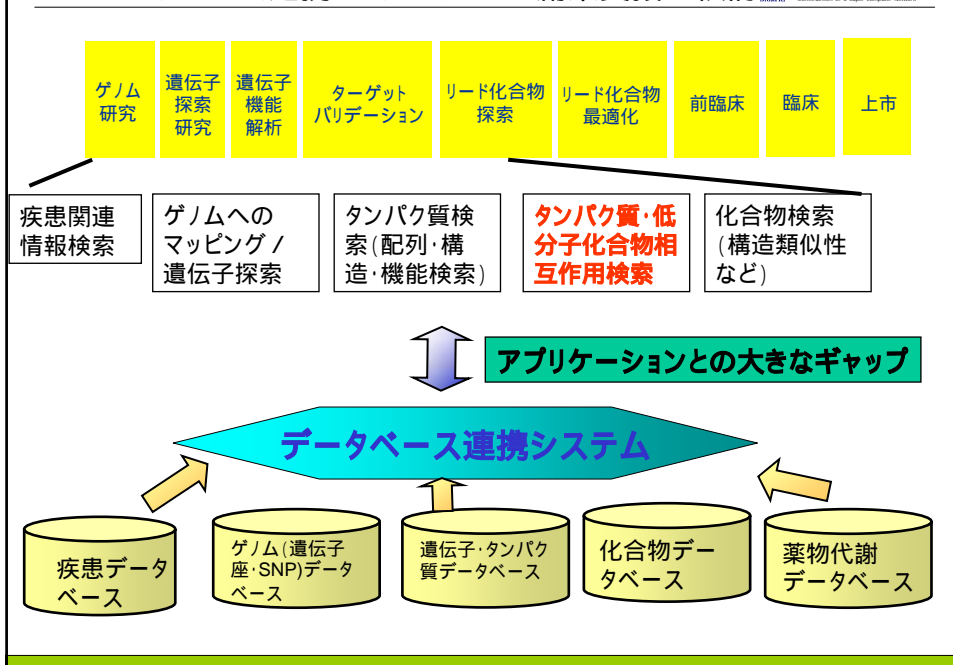
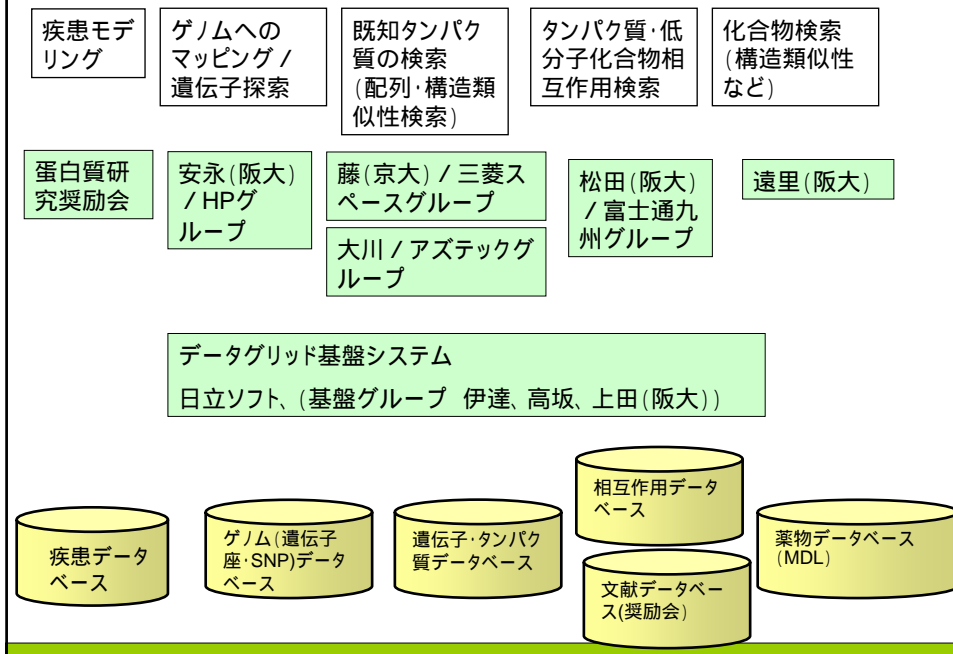


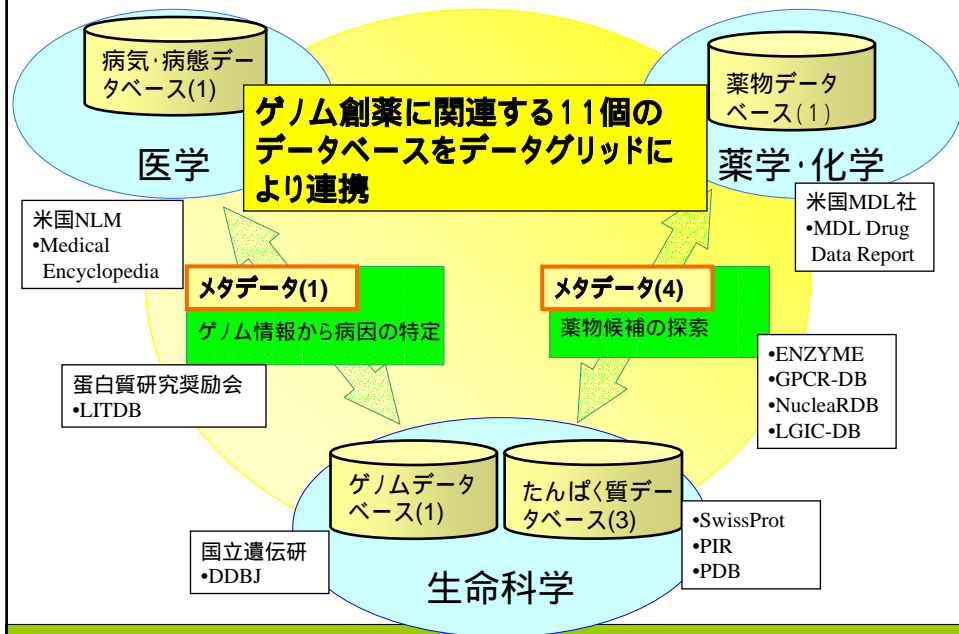
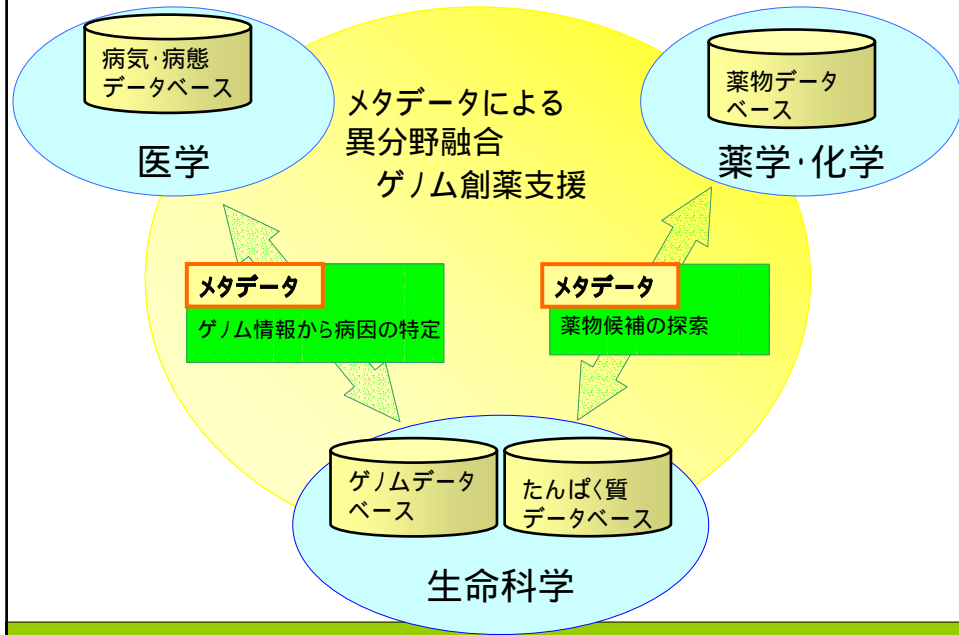
SNP(単一塩基多型)

Genomic Data							Transcription Data					
SNP ID	Coord. Accession	Pos in Contig	Str	5' Flanking Sequence	3' Flanking Sequence	Validation	DNA Chg	AA Chg	Type	mRNA Accession	Protein Accession	Pos in Prot
rs1794906	NT_002933.10	25722157	-	GAGGCTGCTC	GAGGCGCGC	by-substitit	G/C	-	intron	NM_000665	NP_000556	-
rs17636	NT_002933.10	25723676	-	TGGGGGTGCC	CAGGCTACG	by-frequency	C/T	-	coding-synonim	NM_015811	NP_058386	2 592
									coding-antisense	NM_000665	NP_000556	3 477

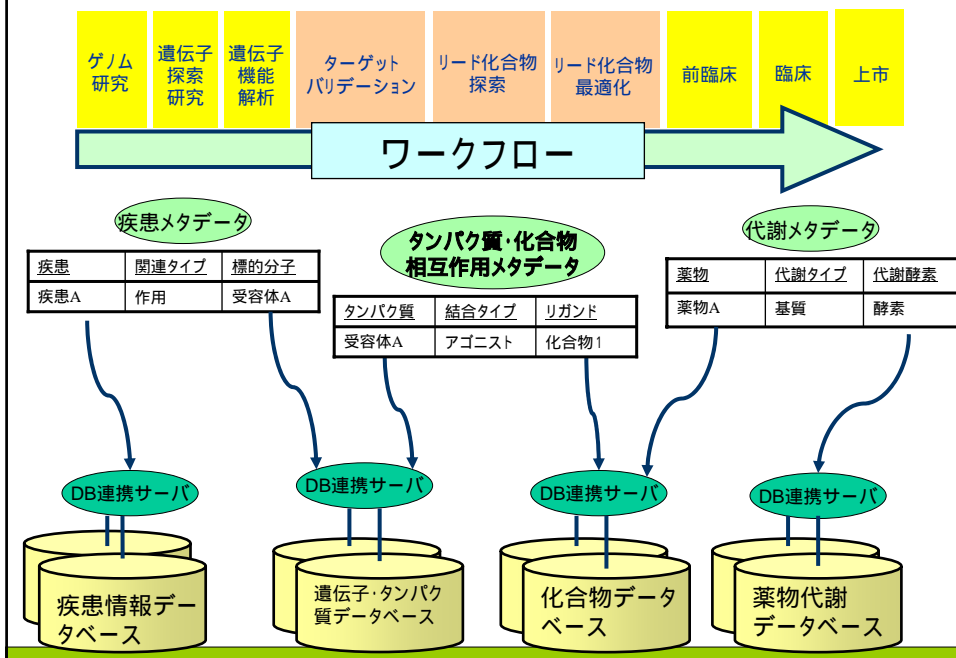
ゲノム創薬に必要なデータベース



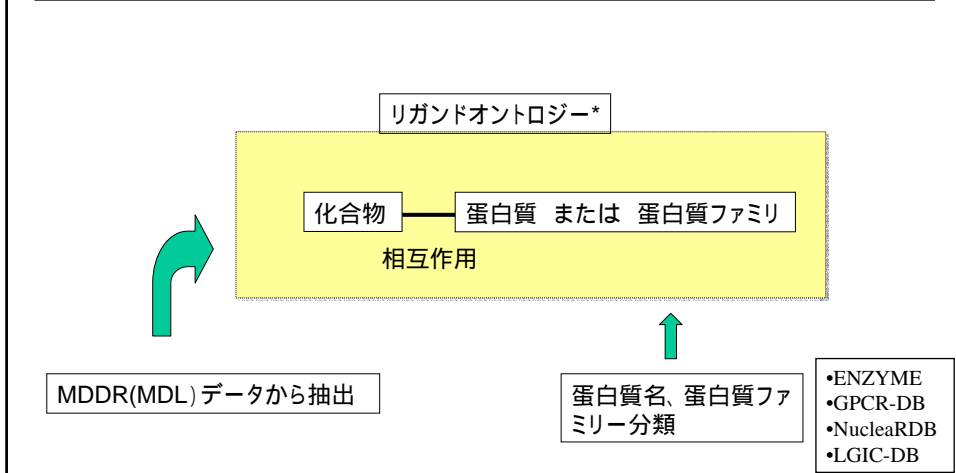




2段階のデータベース連携

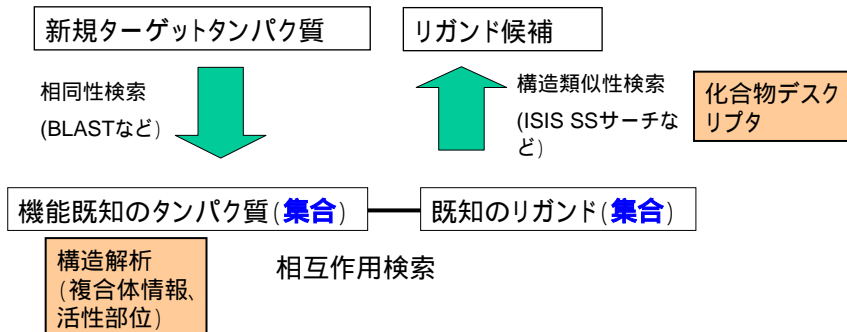


タンパク質・化合物相互作用の表現



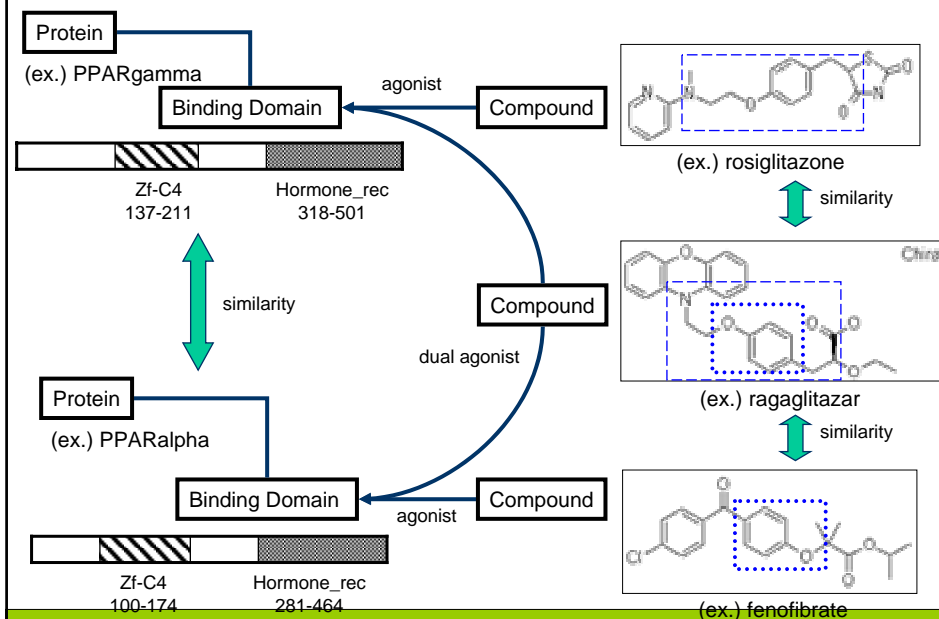
* Schuffenhauer A, Zimmermann J, Stoop R, van der Vyver JJ, Lecchini S, Jacoby E. "An ontology for pharmaceutical ligands and its application for in silico screening and library design," *J Chem Inf Comput Sci*. 2002 Jul-Aug;42(4):947-55.

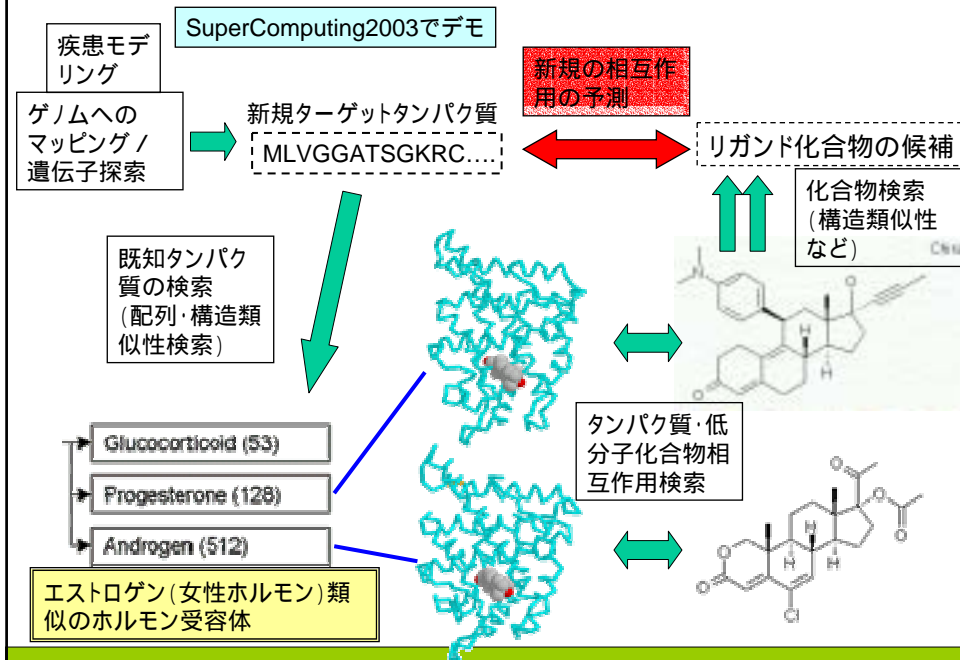
活性相同性 (Activity Homology)



Schuffenhauer A, Floersheim P, Acklin P, Jacoby E.,
 "Similarity metrics for ligands reflecting the similarity of the target proteins",
J Chem Inf Comput Sci. 2003 Mar-Apr;43(2):391-405.

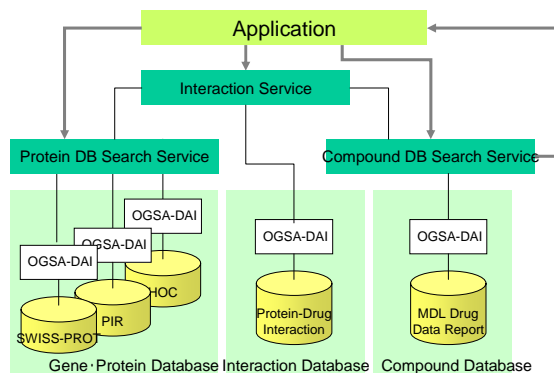
活性相同性の例





Demonstration System

- Bio-related databases are classified into categories.
- Querying databases for each category are provided as our [Grid Service](#).
- The Grid Services are integrated by using **Globus Toolkit 3** with **OGSA-DAI** (Data Access and Integration)



Category	Database	Amount
Disease	Medical Encyclopedia	3079 entries
Genome	DDBJ	Human 7037852 entries, 10176023644 bases
		Mouse 5063486 entries, 6071844270 bases
Protein	Swiss-Prot	137885 entries, 50735179 amino acids
	PIR	283227 entries, 96134583 amino acids
	PDB	23073 entries
Compound	MDL Drug Data Report (MDDR-3D) Ver. 2003.2	142553 entries
Interaction	Ligand Ontology	ENZYME, GPCR-DB, NucleaRDB, LGIC-DB



BioDataGrid



Please come to see our demonstration.

SC2003デモシステム トップページ

BioDataGrid

Protein-Compound Interaction Search

Overview

The BioDataGrid provides a cooperative search among molecular biology databases. This system is compliant to the OGSA which is a standard architecture of grid technologies. You need not to be aware of location or heterogeneity of databases.

Protein-Compound Interaction Search is an application on the BioDataGrid, to find interactions between proteins and compounds from protein view, disease view, or compound view.

Available Databases

Category	Database	Amount
Disease	Medical Encyclopedia	3079 entries
Genome	DDBJ	Human 7037852 entries, 10176023644 bases
		Mouse 5063486 entries, 6071844270 bases
Protein	Swiss-Prot	137885 entries, 50735179 amino acids
	PIR	283227 entries, 96134583 amino acids
	PDB	23073 entries
Compound	MDL Drug Data Report (MDDR-3D)	142553 entries

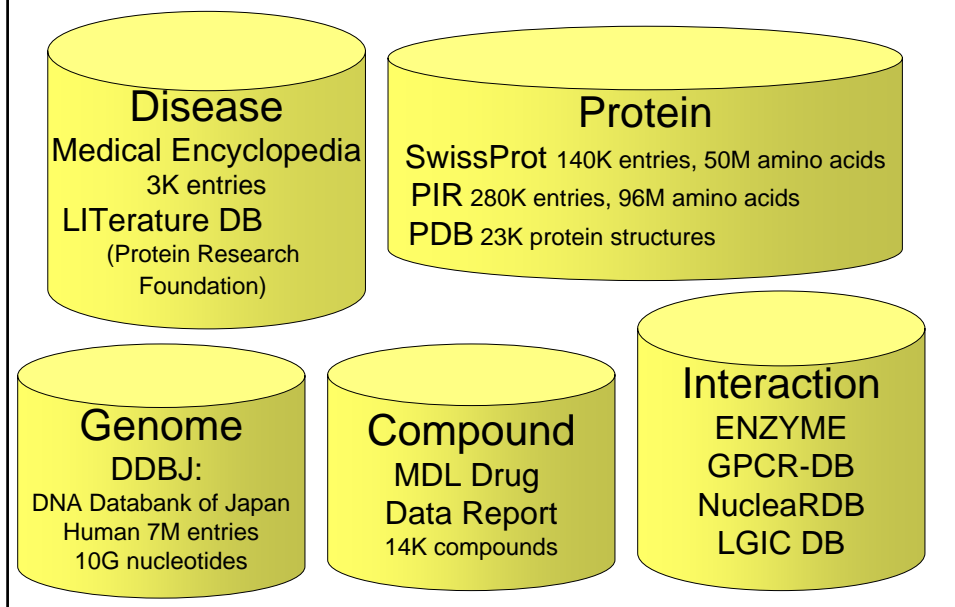
Current Keyword

SWISSPROT ID

Disease Name

デモシステムで利用したデータベース

21



デモシステムのソフトウェア構成

22

- Disease View:** Select a disease. Retrieve a protein (**target protein**) related to the disease.
- Genome View:** Check its genome location.
(adjacently-located genes may also be functionally- related)
- Homology Search View:** Search for similar proteins against DB.
- Protein-Compound Interaction View:** Extract compounds bound to the proteins (using protein-compound interaction metadata).
- Compound Search:** Search for new compounds possibly-interacted to the the target protein.

- Web上に散在する多数のデータベースを連携するデータグリッド技術を開発。
- ゲノム創薬支援を具体例に、疾患、ゲノム、タンパク質、薬物に関連した11個のデータベースを実際に連携することにより、個別のデータベースを意識させることなく相互のデータを動的に関連付けて検索できるシステムを開発し、SuperComputing 2003でデモを行った。
- 最新のグリッド技術であるGlobus Toolkit 3/OGSA-DAIをいち早く取り入れ、実用的な応用システムを構築することにより、データベース連携のためのデータグリッド技術の有効性を実証。

今後の課題

- 応用面：
 - メタデータの作成支援 (統一的な概念記述、オントロジ)
 - ゲノム創薬の探索過程での統一的なスコアリングによる絞込み (優先度付き探索機構の実現)
- グリッド基盤技術：
 - セキュリティ技術の確立
 - 非同期検索(Notification機能)の利用
 - XML DBMS技術の確立