

# 生命科学系データベース横断検索サービスの現状と今後

○大波純一<sup>1</sup>、杉崎太一朗<sup>2</sup>、青木健一<sup>2</sup>、平井信一<sup>2</sup>、牧口大旭<sup>2</sup>、奥村利幸<sup>2</sup>、川本祥子<sup>3</sup>、畠中秀樹<sup>1</sup>、三橋信孝<sup>1</sup>

1.独立行政法人科学技術振興機構バイオサイエンスデータベースセンター(NBDC)、2.三井情報株式会社、  
3.大学共同利用法人情報・システム研究機構ライフサイエンス統合データベースセンター(DBCLS)

## 現在の生命科学データベース横断検索システム概観

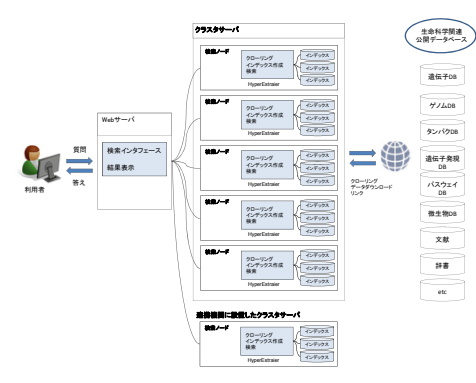
### 検索インターフェース

検索インターフェースのスクリーンショット。日本語検索、ディープWeb検索、お気に入り検索、DBリストダウンロードなどの機能が示されています。

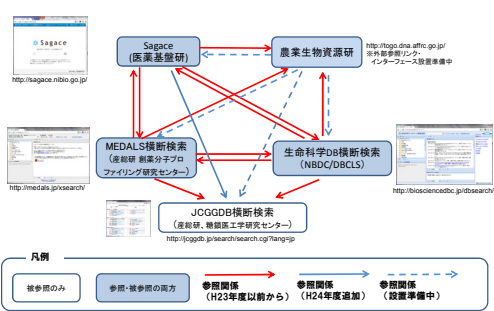
### 機能と特徴 (赤字:機能更新・追加)

機能と特徴の概要。日本語検索、検索語サジェスト、検索結果の絞込み、P2Pノード分散配置、対象DB:378件、公開DBサイト監視、利用マニュアル設置などが示されています。

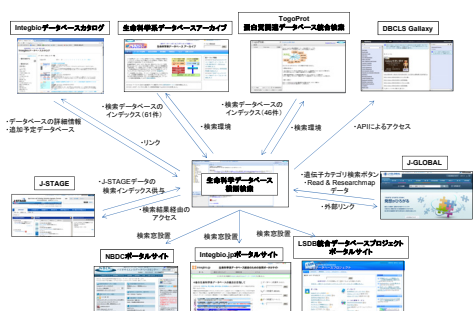
### システム構成



### クラスターサーバの組織間での相互乗り入れ図



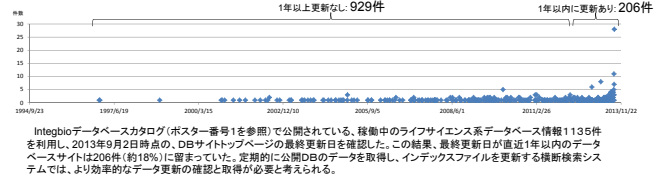
### 外部サービスとの連携



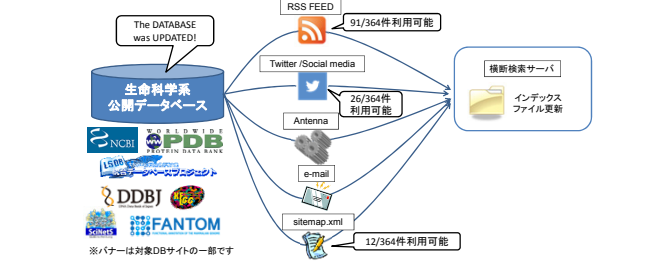
統合的な検索に向け  
1. 幅広い分野のデータベースの相互接続  
2. メタ情報統合化の基準を統合DBプロジェクト内で検討

## これからのシステム機能強化方針

### データベースの更新頻度分析



### インデックスファイル更新の効率化/自動化案



### そして理想的な統合検索へ...

【現状】人間によるデータベースの取得範囲の調査・取得作業

無構造なデータベースでも直接アクセスし、要素を見分ける

取得するデータの例 (赤枠: 検索インデックスに登録する範囲)

Entry	142	COX	10100
Gene name	POLAI	POLA1	p180
Definition	polymerase (DNA directed), alpha 1, catalytic subunit (EC:2.7.7.7)		
Orthology	K02320	DNA polymerase alpha subunit A [E:0.2.7.7]	
Organism	tax	Homo sapiens (human)	
Pathway	tax	M00230	Purine metabolism
Module	tax	M00281	DNA polymerase alpha / prime complex
Class	Metabolism: Nucleotide metabolism: Purine metabolism	[PATH:hs00230]	
Other DBs	NCBI: GI	108507301	NCBI: GeneID: 5422
	ENR: D13240	NCBI: 9113	
	IPRD	02418	Ensembl: ENSG00000101888
	RefSeq	OT11618	NCBI: NM_000002127   UniProt: P08884
Position	chr1	10113	
AA seq	M	APFVQDQSLSDSGEYVSRARPKKSKYGRDELRLKAKAGEKYYKVEEDTGVYEE	
	V	DEEYSLVAVRQDDQNNVDVGGVYDGRDFDDLDLADLADKDGKGRKRMK	
	R	PKVSLAVYRFRNSMPTACAGKDKTADRAVRLSDGLLGDLDLINTETPLGTPHPHNE	
	K	AKKRSRSPSPFVWITAFVGGQAKAVSRKSPKSTVYVAKRSEADKQVVEVTE	
	L	EGEGAMEFEGDQDFEMVEVLEFLMAAKAKRKESEPAEYVNGEADSKQIVSYLSSG	
	L	LVQVCHKDDKGGDQSPFVWVQVDSLPLVYKAGEKQHPHFLVLAHKKDQVQKQV	
	L	LFQVGSIAETHVSCVMKNEKRLVFLPMEKDLATGKETITRPMKQVYVEDEK	
	W	IAVYVMSKSPVENEHAFSPDQVPESENLVWVQEMPLDQDLDETSHPVGTIN	

現在横断検索でも、人手で最新の情報を確認し、公開データベースのエントリーごとに重要な要素を見分けて検索インデックスに登録している。このため、一般的なWeb検索と比べ、高い網羅性と精度で生命科学系データにアクセスすることができる。但しこの手法では、統合的な検索(DB内項目の意味を理解した検索方法)とならない上、情報確認に時間がかかる。

【将来】データベースの構造化・公開データベースごとのRDFメタ情報スキーマの埋め込み

生命科学系データの統合的な利用  
ユニバーサルで負担の少ない取得機構

RDF化・構造化された公開データベース

今後は、統合的な検索の実現のため、構造化(RDF化など)されたデータベースへの取得方法の検討や、生命科学分野の研究に有用な統合項目の検討、データベース取得機構のユニバーサル化などを実施し、情報確認の効率化を推進する。また、連携機関、検索対象データベース管理者、利用者からの意見も積極的に取り入れて、更なる改善に努める方針である。

### robots.txtを公開しているDBサイトの割合

