

Pascal RIGOLET

Bat 112 Institut Curie – Paris-Saclay

06 20 65 77 70

pascal.rigolet@u-psud.fr

BADE

(**BA**se de **D**onnées pour l'**E**nvironnement)

Apprendre à créer et gérer une Base de Données

⇒ **Présentiel + Distanciel (+ Comodal?)**

Mesures Covid19 à respecter

Connexion **windows**

User : prénom.nom@u-psud.fr

(ou : prénom.nom@universite-paris-saclay.fr)

Password : celui de votre messagerie

BASES ET BANQUES DE DONNEES

Objectifs de l'enseignement

- Apprendre à créer et gérer une Base de Données
- Nous disposons de 8 séances pour tout faire
 - ⇒ **L'ACCENT SERA MIS SUR LA PRATIQUE**
- Le site web de l'enseignement sera bientôt opérationnel
- **Contrôle des connaissances** : vous rendrez **un projet** en 3 parties
 - 1 - Création d'une **base de données** personnelle [coef 3]
(application : Environnement/Pollution Chimique/Géologie/Ecologie)
 - 2 – **Rapport de 12 pages** sur cette base de données [coef 2]
(*objectifs ; de la création à l'utilisation ; format pdf*)
 - 3 – **Oral de présentation du projet** [coef 2]
(*dont démonstration de l'utilisation de la base de données + diaporama*)
- On commencera les projets en classe

Système de Gestion de Base de Données (SGBD)

A – Introduction

- Qu'est-ce qu'on entend généralement par **BASE DE DONNEES** ?
- Petit historique
- Limites, fiabilité et stockage des données
- Vers la structuration des données
- Quelques exemple de bases connues
- Une base de données est généralement associée à des services
- Quelques mots de vocabulaire
- Codage ASCII

Base de Données



The collage features several elements: a topographic map with elevation contours; a protein structure visualization with a red and yellow core; a black USB drive; a blue credit card with 'BANQUE CRYPTO' and 'CARTE BLEUE' branding; and a perspective view of a server room with blue lighting.

The interface includes a sidebar with 'Lieux ponctuels', 'Lieux linéaires', 'Lieux surfaciques', 'Zones marines Quadrige', 'Massees sur DCE', and 'Zones conchylicoles'. The main area shows a map with red markers and a 'Table des matières' section.

The PDB entry for 2Z9G includes a primary citation: 'Structural basis of mercury- and zinc-conjugated complexes as SARS-CoV 3C-like protease inhibitors.' by Lee, C.C., et al. (2007). The structure is described as a dimeric assembly with a molecular weight of 24155.60 and a length of 306 amino acids.

The visualization shows a dimeric assembly of the SARS-CoV 3C-like protease, with two subunits in different colors (green and red) and a central active site. The structure is labeled '2Z9G' and includes a 'Biological Assembly' dropdown menu.

Big Data : tous les deux jours, l'humanité produit autant d'information que ce qu'elle a généré depuis l'aube de la civilisation jusqu'en 2003. Plus de 90% des données disponibles ont été produites ces 2 dernières années; ce volume double tous les deux ans d'où l'enjeu considérable que constitue le **traitement**, **l'analyse**, le **stockage** et le **décryptage** de ces « mégadonnées » (« données massives »)

Chaque année, l'Humanité produit plus de 2000 **Exabit** d'informations nouvelles (10^{21} bits, 1000 milliards de milliards de bits), soit 2 Zettaoctets et a déjà stocké **en 2020 plusieurs dizaines de Zettaoctets** d'information (10^{22} bits, 10000 milliards de milliards), environ 40 milliards de téraoctets... Cela a un coût **⇒ Matière + Energie + Eau**

PDB



168599 Biological Macromolecular Structures Enabling Breakthroughs in Research and Education



Advanced Search | Browse Annotations



- Welcome
- Deposit
- Search
- Visualize
- Analyze
- Download
- Learn

A Structural View of Biology

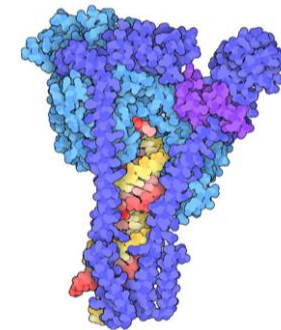
This resource is powered by the Protein Data Bank archive-information about the 3D shapes of proteins, nucleic acids, and complex assemblies that helps students and researchers understand all aspects of biomedicine and agriculture, from protein synthesis to health and disease.

As a member of the wwPDB, the RCSB PDB curates and annotates PDB data.

The RCSB PDB builds upon the data by creating tools and resources for research and education in molecular biology, structural biology, computational biology, and beyond.



September Molecule of the Month



SARS-CoV-2 RNA-dependent RNA Polymerase

Latest Entries

As of Tue Sep 08 2020



Features & Highlights



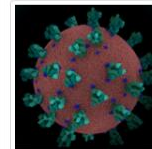
Explore Sequence-Structure Relationships
Protein Feature View maps PDB structures and corresponding UniProtKB sequences and other annotations in an updated interface



Real Time Structural Search of the PDB
A new paper describes a novel method

News

Publications



Virtual Boot Camp
A week-long look at COVID-19 evolution and structural biology has been described in *Biochemistry and Molecular Biology Education* » 09/08/2020

Congratulations, Sophia! » 09/01/2020

Systeme de Gestion de Base de Données (SGBD)

E –Les données dans la base

- ▶ Quantité et qualité des données
- ▶ Performance de la base
- ▶ Exploitabilité des données
- ▶ Le moindre détail peut avoir de l'importance
- ▶ Meilleur compromis

Qu'est-ce qu'une base de données?

Système de Gestion de Base de Données (SGBD)

E – Quelques définitions

- Une **base de données** c'est essentiellement une **collection structurée d'informations** non nécessairement du même type (format)
- Une base de données est usuellement localisée en un seul lieu et un seul **support** (dupliqué en fait) qui est généralement informatique (numérique)
- Pièce centrale des dispositifs informatique qui servent à la collecte, au stockage et à l'utilisation des informations
- **SGBDR** : logiciel moteur qui pilote la base et en permet la manipulation et l'exploitation (interrogation)
- Toutes les secondes le volume d'information ne cesse d'augmenter contribuant au **Big Data**. Mais sans analyse et sans base de données (contribuant à préparer les données à leur analyse), le Big Data n'est rien.

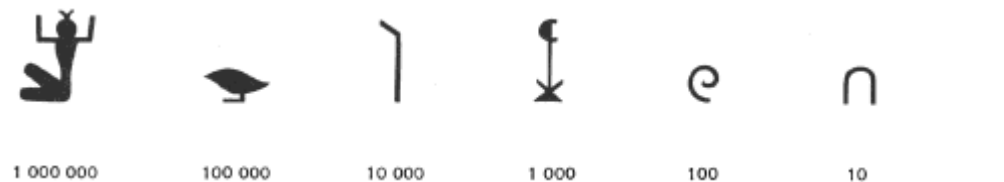
Systeme de Gestion de Base de Données (SGBD)

A-2 – Quelques exemples de bases de données identifiables

- Décomptes et registres (Babylone, Haute Egypte, antiquité, moyen âge)



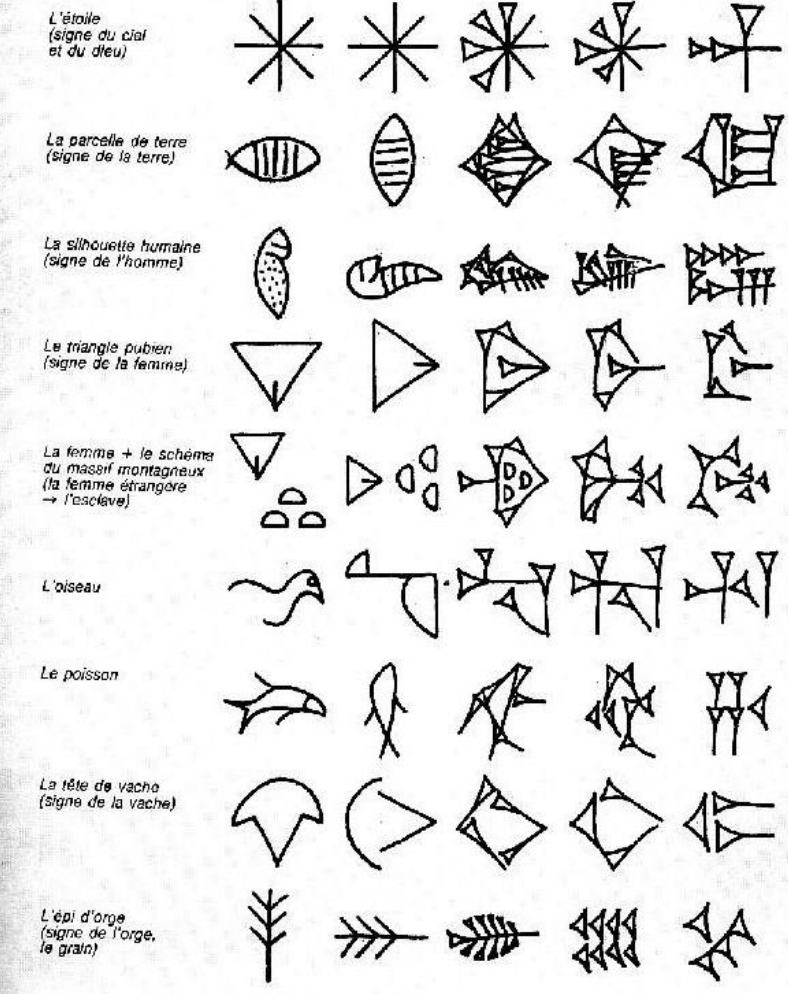
Babylone, écriture cunéiforme



Les chiffres et nombres, haute Egypte



Dates approximatives : -3300 -2800 -2400 -1800 -700



1	7	11	<7	21	<<7	31	<<<7	41	<<<<7	51	<<<<<7
2	11	12	<11	22	<<11	32	<<<11	42	<<<<11	52	<<<<<11
3	111	13	<111	23	<<111	33	<<<111	43	<<<<111	53	<<<<<111
4	1111	14	<1111	24	<<1111	34	<<<1111	44	<<<<1111	54	<<<<<1111
5	11111	15	<11111	25	<<11111	35	<<<11111	45	<<<<11111	55	<<<<<11111
6	111111	16	<111111	26	<<111111	36	<<<111111	46	<<<<111111	56	<<<<<111111
7	1111111	17	<1111111	27	<<1111111	37	<<<1111111	47	<<<<1111111	57	<<<<<1111111
8	11111111	18	<11111111	28	<<11111111	38	<<<11111111	48	<<<<11111111	58	<<<<<11111111
9	111111111	19	<111111111	29	<<111111111	39	<<<111111111	49	<<<<111111111	59	<<<<<111111111
10	<	20	<<	30	<<<	40	<<<<	50	<<<<<		

Systeme de Gestion de Base de Données (SGBD)

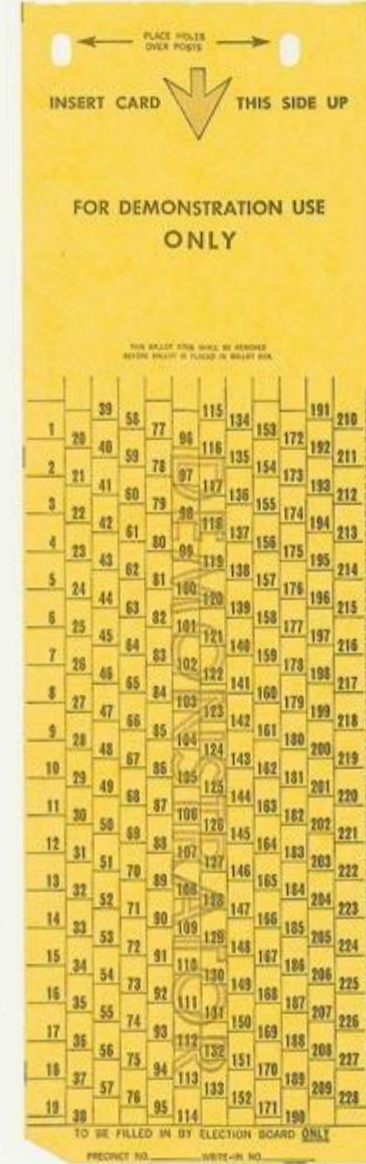
A-2 – Quelques exemples de bases de données identifiables

- Décomptes et registres (Babylone, Haute Egypte, antiquité, moyen âge)
- Archives dans l'histoire
- Coordonnées bancaires (dès le 19^{ème} siècle)
- Compagnies de téléphone
- Clients/fournisseurs
- Horaires des train et réservations de billets
- Fiches de police – Carte d'identité, Passeport, ...
- Sécurité Sociale (vers 1946 avec fiches support papier, fichiers boîtes) puis ordinateurs et fichiers (N° INSEE)
- Premier ordinateurs commandés par les USA pour applications civiles : UNIVAC recensement de 1950
- Loi informatique et liberté de 1976
- Google, Yahoo et autres « moteurs de recherches »
- Face book (et ses dangers)
- Apogée et tous les fichiers gérant l'inscription et les notes des étudiants
- Edwige (un projet avorté?)
- **Chaque français en moyenne est enregistré dans une centaine de DB**

Système de Gestion de Base de Données (SGBD)

A-2 – Quelques exemples de bases de données identifiables

L'Univac (1951)



Gogol (« Googol ») : 10^{100} \Rightarrow grand nombre versus infini

- ▶ Très supérieur au nombre de particules dans l'Univers connu (environ 10^{80}) et pourtant simple à écrire (1 suivi de 100 zéros)
mais très inférieur au nombre de cubes de Planck (côté de la longueur de Planck : $1,616 \times 10^{-35}$ m) dans l'univers observable.
- ▶ Mathématicien Edward Kasner (« *Mathematics and the Imagination* » - 1938)
- ▶ Un peu de légende geek (son neveu de 9 ans aurait trouvé ce nom)
- ▶ Origine du nom explicitement revendiqué par les fondateurs de Google
- ▶ 10^{gogol} : « gogolplex »
- ▶ Une certaine écriture de l'infini?...
- ▶ Si unité = mètre : infiniment petit (connu) : 10^{-32} ; infiniment grand (connu) : 10^{32}
- ▶ On ne peut pas représenter l'infini dans une base de donnée : monde concret/connu
- ▶ Génome humain : 23000 gènes ; 3,2 milliards de pdb ; 750 Mo (haploïde)
- ▶ Diamètre de l'univers observable : 9.2×10^{10} années-lumière

Systeme de Gestion de Base de Données (SGBD)

E – Fiabilité des données

Hélas, les données ne sont ni fiables, ni éternelles...

- ▶ Utilité d'une base de données (il faut qu'elle serve à quelque chose)
- ▶ Sécurité des données
- ▶ Contenu = vérité ?
- ▶ Stockage fichiers en plusieurs exemplaires (non nécessairement identiques)
- ▶ Versions ; date mise à jour
- ▶ Stockage : altération avec le temps
- ▶ Autres problèmes ...

Systeme de Gestion de Base de Données (SGBD)

E – Fiabilité des données

Support de l'information : un problème à moyen et long terme...

► **Conservation des données**

Papier : ~ 1000 ans

Pierre : ~ 10000 ans

Bambou & papyrus : ~ plusieurs milliers d'années

Peaux : plusieurs centaines d'années

Grottes : > 35000 ans (Chauvet)

Toiles, peintures, mosaïques : 1000 à 3000 ans

Argile (écriture cunéiforme) : > 4500 ans

Film, photos, ciné : ~ 100 ans (demandent ensuite restauration)

CD : 5 à 10 ans maxi

Clé USB : 0 à 20 ans maxi

Disquette : 15 ans

Bandes magnétiques : 100 ans et 1000 ans (nouvelle génération)

Disque dur : 5 à 10 maxi

Cloud : ??? (dépendant du reste; supports changés régulièrement)

ADN : 300000 ans

Quartz (gravure) : > 1 milliard d'années!!

- Produire, stocker et exploiter **les data consumeront** d'ici un siècle à un siècle et demi environ toute la matière et toutes les **ressources de la planète !!!**

Système de Gestion de Base de Données (SGBD)

E – Quelques définitions

Une **base de données (relationnelle)** c'est essentiellement une **collection (structurée) d'informations** non nécessairement du même type (format)

(Il est vrai qu'il existe des bases de données non structurées pour le big data mais qu'il faut passer par une structuration des données pour pouvoir les traiter, les analyser et les exploiter)

⇒ Il va donc falloir **apprendre à stoker et organiser l'information** en vue de son exploitation : il existe des méthodes pour cela; cela ne s'improvise pas.

⇒ **Pas de place pour le bazar...**

⇒ **mais la logique et le bon sens ont tous leurs droits!**