# DBW – Databases and Web development



MSc Bioinformatics for Health Sciences

# Aims

- Review a number of technologies to handle bioinformatics data:
  - Computer communication, design of web applications, basic database design and optimization.
  - This is NOT a programming course, it is about designing and building applications in an heterogenous scenario

- The final objective is to built a **fully operative application** using the appropriate combination of the techniques reviewed.

# Bioinformatics & Internet

- Tools and data should be available through web

- Ex. Nucleic Acid Research reviews:
  - Database Issue (January) 1170 DBs
  - Web Server Issue (July) 1200 Servers

# NAR Database issue recommendations

- "The pre-submission enquiry must present a working web accessible database "

- "The quality, quantity and originality of data as well as the quality of the web interface are the most important. Good data with a poor interface or vice versa are never sufficient for consideration. "

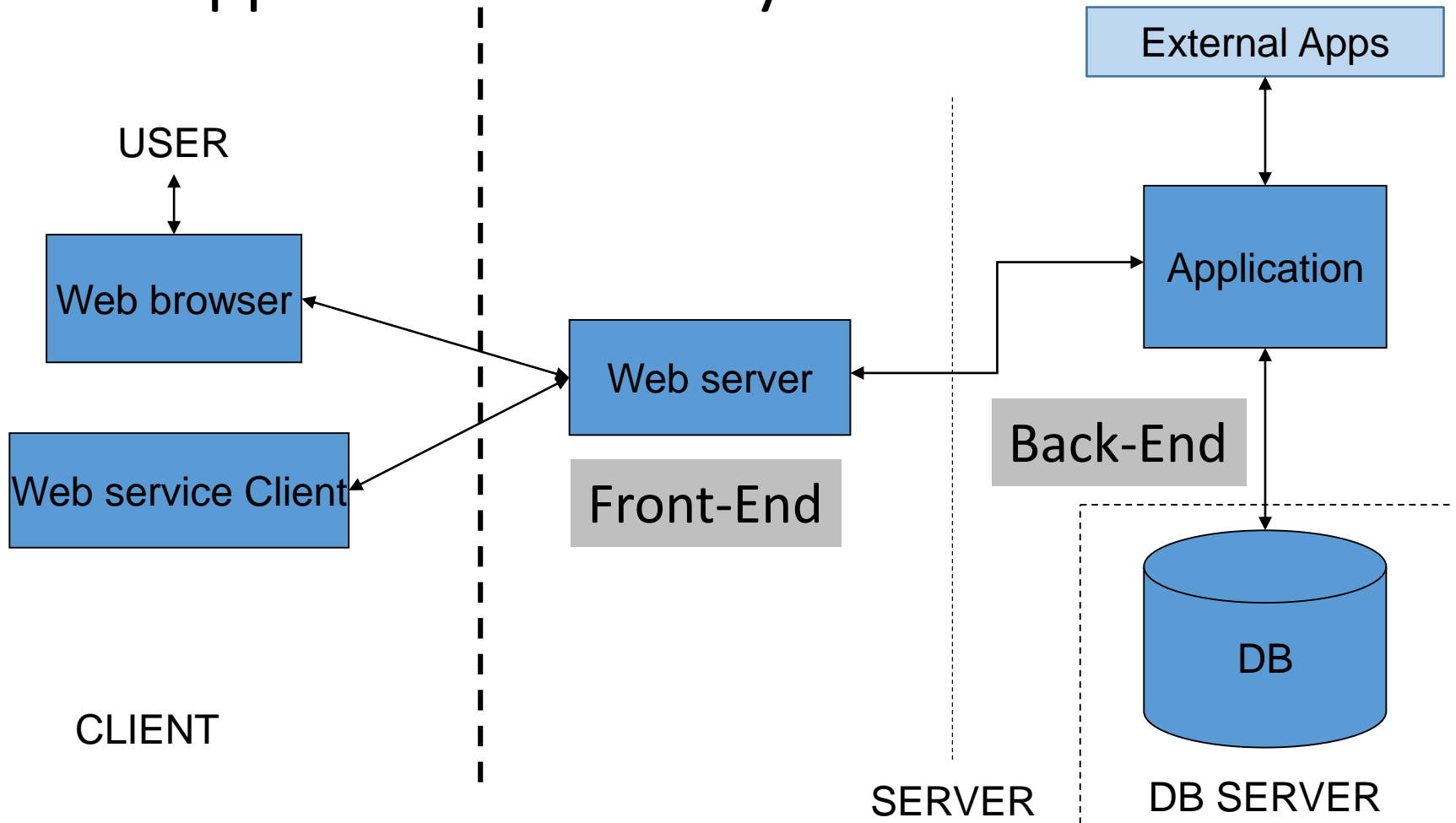- "Do get a domain name for your website. URLs to specific IP addresses/ports are unlikely to stand the test of time."

- (…)

*Nucleic Acids Research*, Volume 35, Issue suppl_1, 1 January 2007, Pages D1–D2

https://academic.oup.com/nar/article/35/suppl_1/D1/1088333

# Web applications by access type

- Web interfaces
  - Provide a user friendly interface (web based) to "human" users
    - Users known how to use the interface
    - There is no need to install software
    - Single operations (no large scale)
    - Must adapt to navigation uses (low latency, synchronous answers,…)

- Web services
  - Provide a programmatic interface (using Web protocols)
  - Intented to interact with software, not humans
    - Well-defined data formats required.
    - Adequated for large scale operations

- Modern applications will normally offer both

# Web application anatomy

USER

Web browser

Web service Client

Web server

Front-End

CLIENT

External Apps

Application

Back-End

DB

SERVER

DB SERVER

# Web application styles

- Access to data
  - Friendly interface to data repositories

- Web Interfaces to stand-alone software
  - Collect input parameters and redirect output

- Workbenches (e.g. Galaxy)

- On-purpose applications & DBs

- Web services (programmatic access)

López-Ferrando *et al.* **Nucl. Ac. Res** 2017

# Web interfaces to apps.

http://mmb.irbbarcelona.org/MDWeb2

Bioinformatics. 2012 28(9):1278-9.
doi: 10.1093/bioinformatics/bts139

https://usegalaxy.eu/

# Special purpose applications & DBs

**IMID-KIT [imidclinica] - Microsoft Internet Explorer**

Archivo   Edición   Ver   Favoritos   Herramientas   Ayuda

Atrás ▾   ✕   Búsqueda   ★ Favoritos   📷 ▾ 📷 ▾ 

Dirección  http://inb.bsc.es/imid_kit/egroupware/imidclinica/index.php   → Ir   Vínculos »

Google [G▾]   ✕ Ir  ☆ Marcadores ▾  Pop-ups permitidos  Corrector ortográfico ▾ »   Configuración ▾

Añadir ...

imidclinica

**imidclinica**

✳ INB°  NB°

Referencia:
Especialidad:
Centro:

**Busqueda de Casos**

· Ref. Caso                 contiene ▾
· Ref. Banco DNA       contiene ▾
· Datos completos centro    ☐
· Datos completos              ☐

· Centro(s)   Servicio de Dermatología, Hospital Universitario Gregorio Marañón (M
              Servicio de Dermatología, Hospital General Universitario de Valencia
              Servicio de Dermatología, Hospital Universitario 12 de Octubre (Madr
              Servicio de Dermatología, Complejo Hospitalario Juan Canalejo (A Co
              Servei de Dermatologia, Hospital de la Santa Creu i Sant Pau (Barcelo

Form

Codigo Banco AD
Datos Centro Co
Datos Completos

**Cuestionario cli**

A. Criterios de In
B. Criterios de Ex
C. Antecedentes
D. Variables Clíni
E. Variables Clíni
F. Enfermedades Asociadas (Diagnosticadas y Documentadas) (0/36)
G. Tratamiento (0/155)
G1. Tratamiento Etanercept y Adalimumab (0/48)
G2. Tratamiento Rituximab y Anakinra (0/49)
H. Actividad de la Enfermedad en el Momento de la Extracción (1/75)

**Cuestionario epidemiologico**

Extracción (0/0)
Actividades (0/0)
Hábitos (0/0)
Cuestionario Dem
Estudios (0/0)
Cuestionario Gene

**Cuestionarios clinicos**

Selecciona los campos a incluir en la busqueda

[Expandir Todos] [Colapsar Todos]      [Seleccio

Selecciona los campos a incluir en la busqueda

[Expandir Todos] [Colapsar Todos]      [Seleccionar] [Limpiar]

⊟ Reumatología
  ⊞ A. Criterios de Inclusión
  ⊞ B. Criterios de Exclusión
  ⊟ C. Antecedentes Familiares
      ⊞ Psoriasis        ◉ si ◯ no ( ☐ Cualquiera)
      ⊟ EII              ◯ si ◯ no ( ☐ Cualquiera)
          ⊟ Crohn        ☑
                         padre
              · Parentesco   madre
                         hermanas/hermanos ▾ ( ☐ Cualquiera)
          ⊞ Colitis Ulcerosa  ☐
      ⊞ AIC              ◯ si ◯ no ( ☐ Cualquiera)
      ⊞ Otras Enfermedades ◉ si ◯ no ( ☐ Cualquiera)
  ⊞ D. Variables Clínicas y Biológicas Articulares
  ⊞ E. Variables Clínicas Extra-Articulares
  ⊞ F. Enfermedades Asociadas (Diagnosticadas y Documentadas)
  ⊞ G. Tratamiento
  ⊞ G1. Tratamiento Etanercept y Adalimumab
  ⊞ G2. Tratamiento Rituximab y Anakinra

✳ INB°

**imidclinica**

**Busqueda de casos**

✳ INB°

| #1 | Reumatología > C. Antecedentes Familiares > Psoriasis: si | 3 | Resultado | Grabar set |
|---|---|---|---|---|
| #2 | Reumatología > C. Antecedentes Familiares > EII > Crohn: SI | 0 | | |
| #3 | Reumatología > C. Antecedentes Famil... si | | | |
| | Operadores posibles: O, Y, NO | | | |

[Grabar búsqueda actual] [Nueva Busqueda]

**imidclinica**

✳ INB°

#4. #1 Y #3:

Num Casos: 2

| ID. | Especialidad | Centro | Datos Cuest. | Id DNA Bank | Datos Epid. | Compl. Centro | Completo |
|---|---|---|---|---|---|---|---|
| 30112345 | Reumatología | Unitat de Recerca de Reumatologia, Hospital Universitari Vall d'Hebron (Barcelona) | 22/492 | | 0/0 | | |
| 3012345 | Reumatología | Unitat de Recerca de Reumatologia, Hospital Universitari Vall d'Hebron (Barcelona) | 25/492 | | 0/0 | | |

Volver a la búsqueda

14

# Bioinformatics web-services and workflows

Input: Protein Id and DNA genomic sequence

Building of Blast Database from DNA seq.

BLAST Search

Run GeneWise to detect gene structure

# Building a (web) application

1. Define specifications
2. Analyze data and built a data model
3. Decide/prepare Database implementation
4. Build ETL if necessary
5. Define interfaces
6. Define and prepare files/scripts layout
7. Write application code
8. Test, debug, document…

# Course logistics

- Web site(s)
  - Course materials:
    - http://mmb.pcb.ub.es/formacio/

  - Personal sites:
    - http://mmb.pcb.ub.es/formacio/~dbwXX

  - SSH Access
    - ssh mmb.pcb.ub.es –p 22021 –l dbwXX

  - MySQL Access
    - Localhost only

# Software to install

- Ideally Linux (may need root privileges)

- From Linux distribution
  - Apache Web Server (v. 2.x)
    - With PHP 5.x and mysql support
  - MYSQL server (v. 5.x)
  - MYSQL Workbench or phpMyAdmin

- Netbeans (PHP module) (optional)

- MongoDB (optional)

# Evaluation

- Exercices, in-class projects (20%)

- Personal web site (20%)

- Web application project (60%)
  - Progress presentations
  - Fully operative web application using DBs

# Evaluation

- Web application project
  - 3-4 people / group
  - Free subject (bioinformatics preferred)
  - Should include DB management, web interface, users management
  - May use fake data if necessary
  - Available at the personal web site
  - Preferred languages: PHP, Perl, Mysql
  - Source code at github

# Evaluation

- Web application project
  - Steps:
    - Initial specification
    - Data analysis & Database design
    - Project Demo
    - Final application
- Installed on server
  - mmb.pcb.ub.es/formacio/~dbwXX
    - Account dbwXX

# Basic computer communication protocols

# Aim & Outline

- Understand the basic components of computer communication protocols
  - Concepts of client and server
  - Addressing servers and data
    - Computer addresses (MAC Address, IP Address)
    - Ports
    - Resource identification
      - DNS
      - URL/URI concept
  - Client/server transactions
    - HTTP protocol

# Present internet

- Huge network of computers using common communication protocols (TCP/IP, HTTP)

- Distributed, no central servers
  - (Well, not really true in bioinformatics)

- Common language: HTML/CSS (XML)

- Content mostly static, but dynamic behaviour is possible through web applications

# Components



Request

Answer

Client                    Network                    Server

- Client and Server logic and physical addresses
- Data
- Data meta-information
  - Nature of data
  - Request (what to do)
  - Applications involved (email, web, etc.)

# How it works: TCP/IP

- Packet switching
  - Packet switching breaks the signal in small fragments ("packets") each of them containing the complete information about source and destination
    - Packets can share a single communication line

  - Users have the idea of a dedicated line but, in fact, it is not. Of course, the bandwidth is limited.

- Computers connected to internet should have addresses
  - MAC Address: Address of the physical interface
  - IP Address: Address of the computer

# IP addresses

- Allow to find destination irrespective of the nature of the network media.

- Each device has a "unique" IP address

- IPv4: 32 bits (4 x 1 byte (0-255) numbers)
  - Max: $2^{32}$ : aprox 4.3 x $10^9$

  - P.ex. 161.116.222.59  (mmb.pcb.ub.es)
  - 4 levels are hierarchical

- Some addresses are reserved, and some networks are "local"

- IPv6: 128 bits (16 bytes). Max: $2^{128}$ (3.4 x $10^{38}$)

# Names *vs* addresses (Domain Name System)

- IP addresses are not easy. Most hosts have also a "name":

    f. ex. www.ncbi.nlm.nih.edu

- Host names have a structure similar to IP addresses: Top domains (.es, .edu, correspond to full class domains and subnets are indicated by prefixes.

    – ub.es (161.116.x.x)

    – bq.ub.es (161.116.154.x)

    – www.bq.ub.es (161.116.154.18)

# How to address applications in a server: Port ids.

- Each host has **one IP address** but has **several ports** for known services

- Ports are 2-Byte numbers.
  - 0-1023 are "**Well known ports**" (Telnet: 23, FTP: 21, HTTP: 80, …).
  - 1024-49151 are "**Registered ports**", usually managed by applications (MySQL: 3306)
  - 49,152-65,535 are "**Dynamic and/or private ports**" freely usable.

- Communication to ports triggers the specific application to deal with the data

- However, different ports from official ones can be used to:
  - Hide applications
  - Have more then one server in the same IP address
  - Hide servers in internal networks.

# URI/URLs

- Resources must be identified in a way that includes all the necessary details:

http://mmb.pcb.ub.es:80/courses/master.htm#top

| protocol | host | port | resource in the host | label |

Optional

Missing parts of the URL are filled by default!!

# Client – server communication

- Most Web Applications use HTTP (hypertext transfer protocol), although sometimes FTP, SMTP

- HTTP is a client-server communication protocol
  - Link between client and server is dynamic
  - Usually limited to a single transaction
  - Requests composed by a query operation and a variable set of headers.
  - Answers: headers + data

# Client – server communication

- Relevant Operations: GET, POST
  - GET: Simple retrieval, all information/parameters included in the URL
    - Simple queries, static information
    - Required to be used as hypertext links
  - POST: Query defines the resource, but input data follows
    - Input data can be of any type (including binaries, whole files) or size (within limits)

- Relevant HTTP headers
  - Content-type (POST): input data format
  - Content-type (Answer): Data MIME type (text/html, image/jpg, …)
  - Location: Redirects browser
  - Set-cookie: Set a "cookie" on users' software.

# Cookies

- Small information tags sent as HTTP headers and stored in the browser side
  - Are associated with a URL, and are sent back to the server whenever that URL is visited within a expiration date

```
Set-Cookie:
  PHPSESSID=bb56ee648aeac6923e3360a7b8284a6f;
  path=/
```

- Useful to "remember" clients, but some people disables them!