

P2P and Grid Combination

(BioMAJ Peer-to-Peer extension)

Anthony ASSI

HPC R&D Engineer

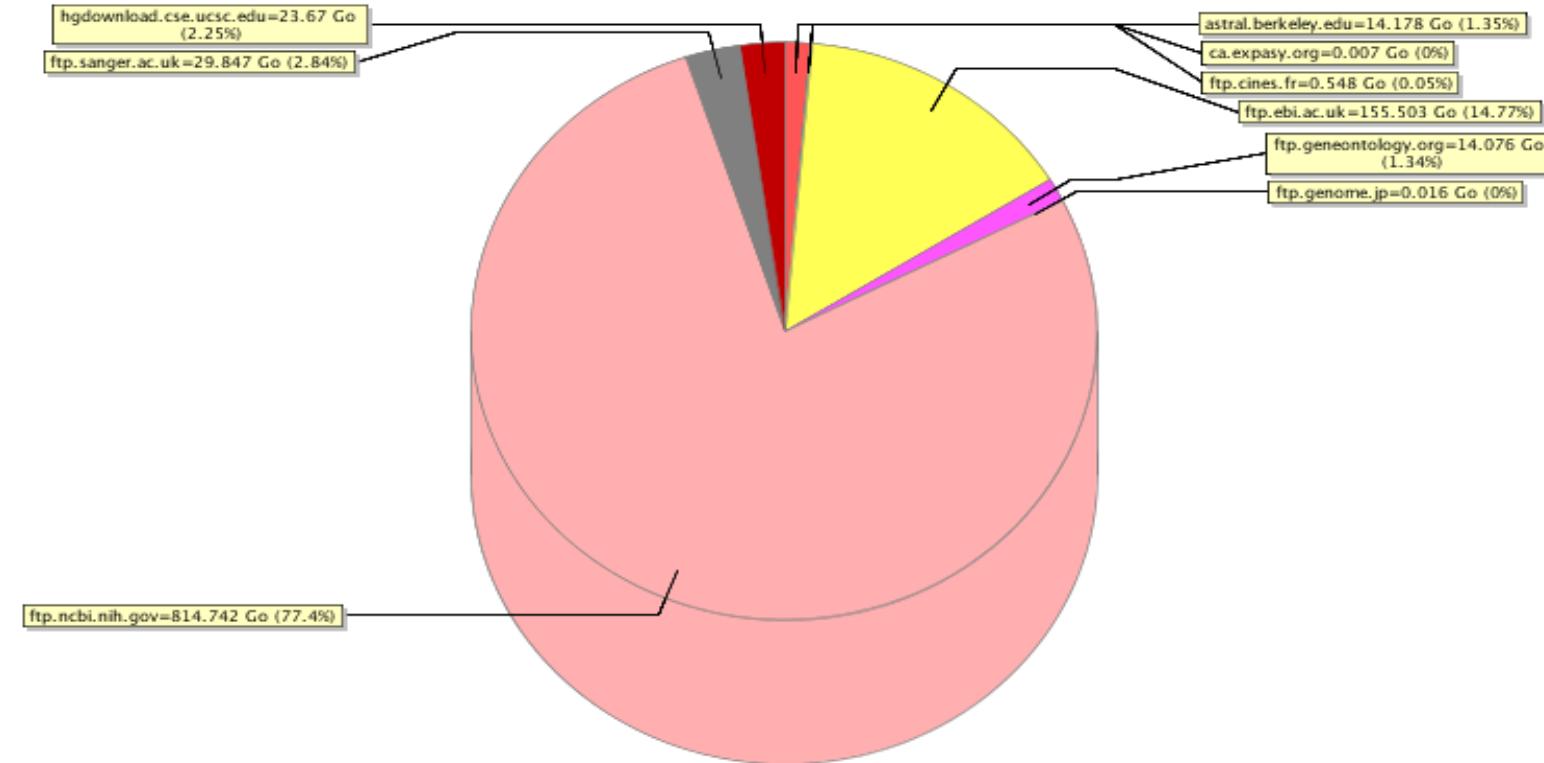
OUEST-Genopole Bio-Informatics Platform
anthony.assi@inria.fr

D. Allouche, A. Assi, Y. Beausse, C. Caron, O. Collin,
O. Filangi, J-M. Larré, L. Legrand, H. Leroy, V. Martin



State of the Art !

Databanks size distribution by Server Name

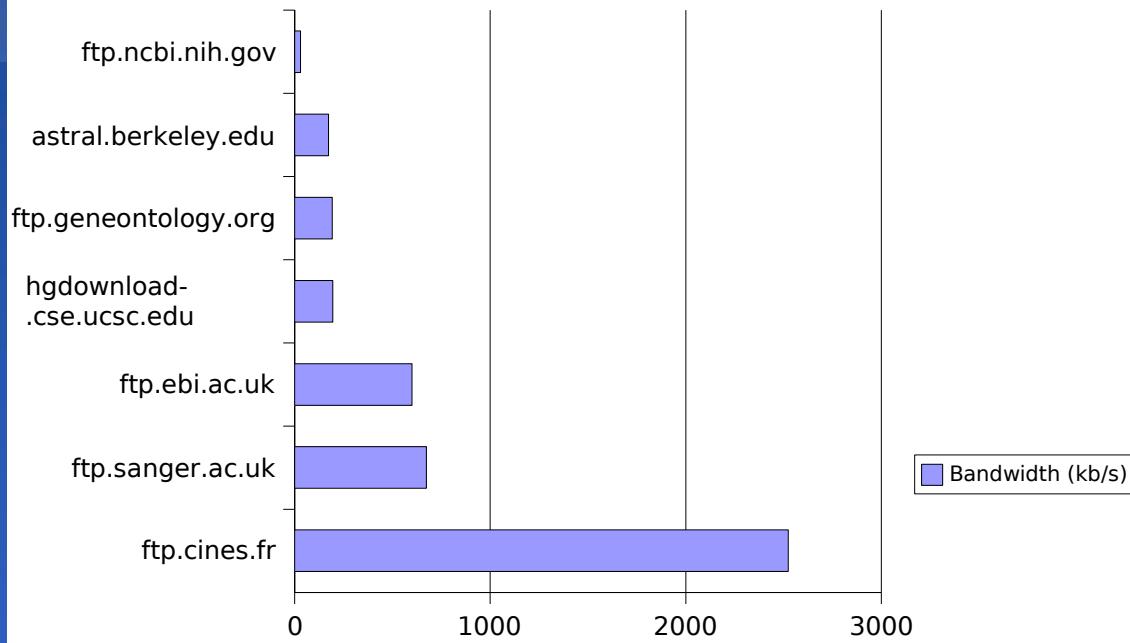


- astral.berkeley.edu
- ca.expasy.org
- ftp.cines.fr
- ftp.ebi.ac.uk
- ftp.geneontology.org
- ftp.genome.jp
- ftp.ncbi.nih.gov
- ftp.ncbi.nih.gov
- ftp.ncbi.nih.gov
- ftp.ncbi.nih.gov
- ftp.ncbi.nih.gov
- ftp.ncbi.nih.gov
- hgdownload.cse.ucsc.edu
- ftp.sanger.ac.uk

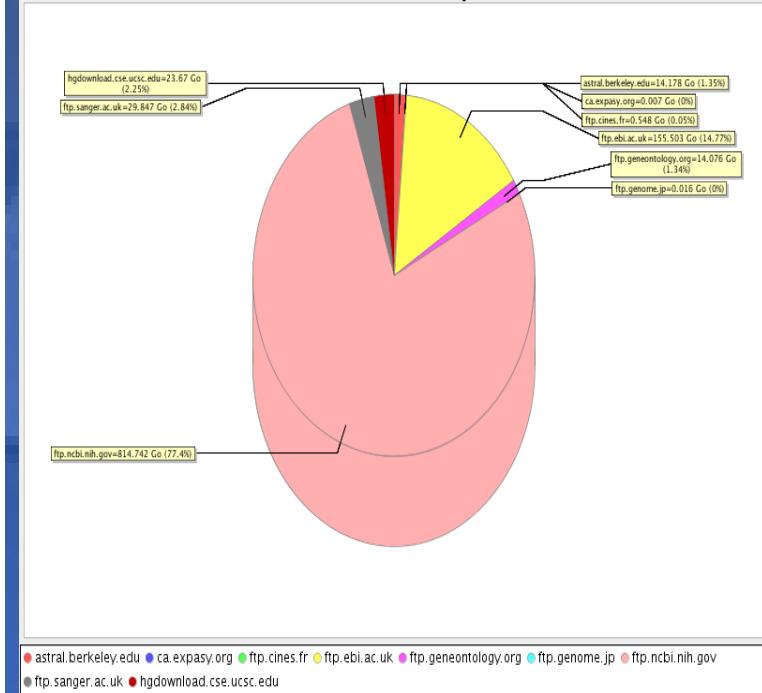
Download Bandwidth Bottleneck

Server name	Bank Name	Bandwidth (kb/s)	Color Code
ftp.ncbi.nih.gov	Genbank, Genomes	31.12	Light Red
astral.berkeley.edu	ASTRAL	173.62	Orange
ftp.geneontology.org	Gene Ontologie	192.76	Magenta
hgdownload.cse.ucsc.edu	GoldenPath	195.22	Brown
ftp.ebi.ac.uk	NR, STS	600.81	Red
ftp.sanger.ac.uk	PFAM	674.12	Grey
ftp.cines.fr	ImMunoGeneTics	2523.59	Brown

Bandwidth Limitation

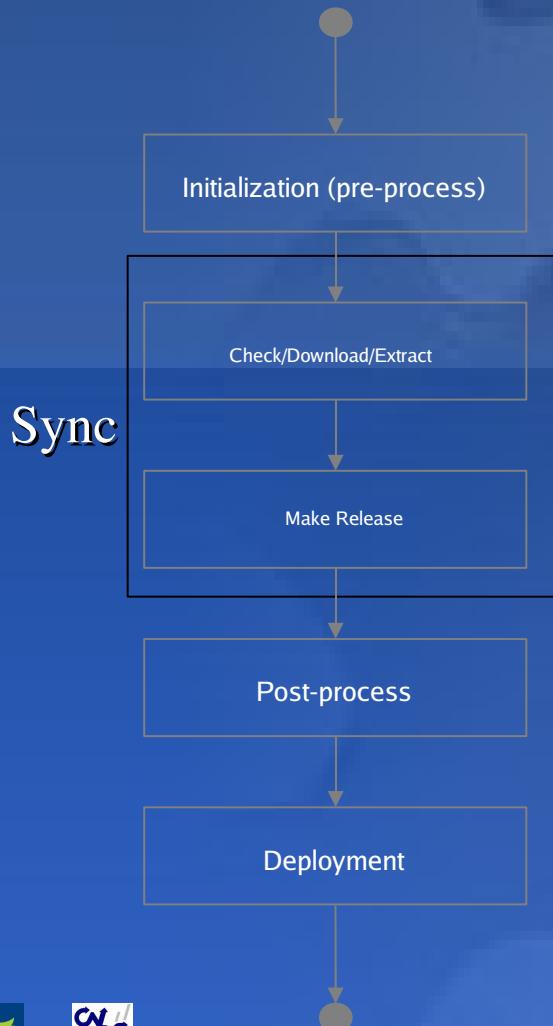


Databanks size distribution by Server Name



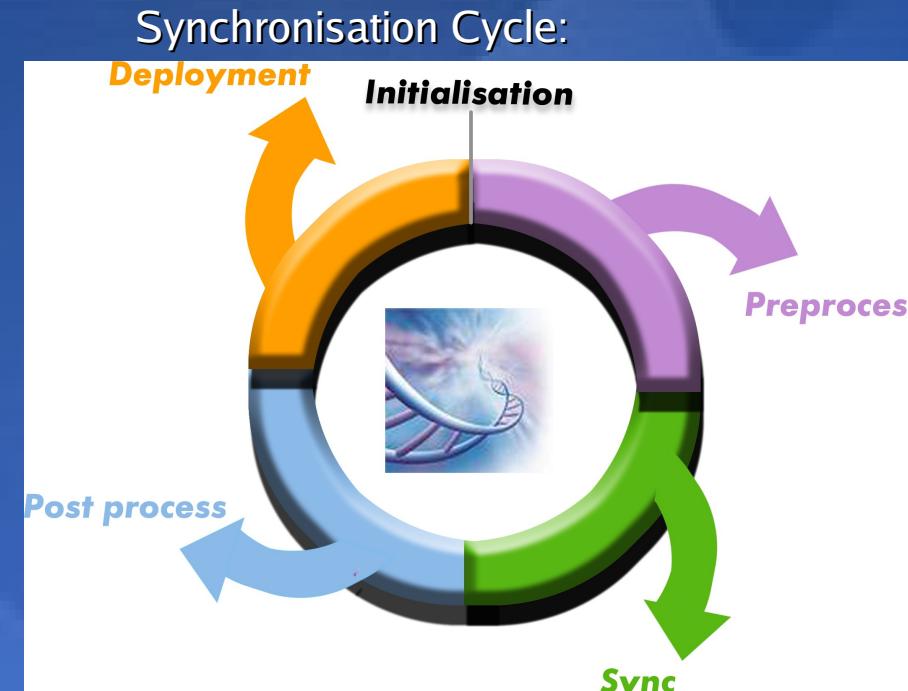
BioMAJ P2P Initiative

The normal WorkFlow:



Goal:

Accelerate the download process of the banks using the Peer-to-Peer technologies to eliminate the Bandwidth Limitation problem among the Web Servers (FTP, HTTP)

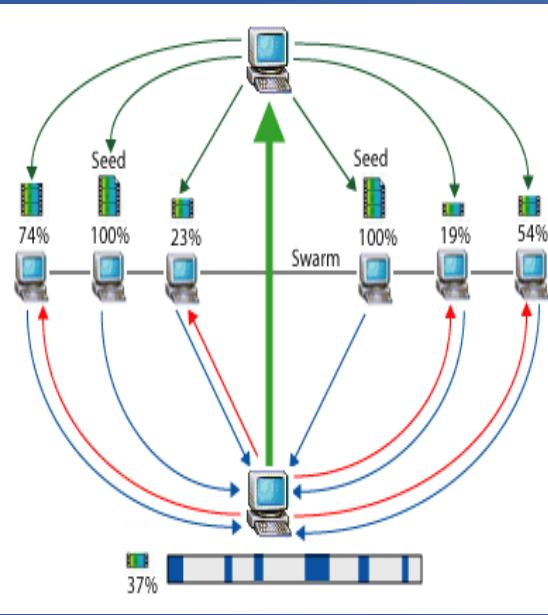


Peer-to-Peer Implementations



Advantages:

Easy to Deploy
Scalable
Give-Give Protocol, thus limiting Free-Riders situations
The network nodes are grouped around one Tracker
Tracker uses HTTP Messages



Disadvantages:

SPOF : Tracker
Open Standard Protocol

The P2P new Workflow



BitTorrent terminology

Peer: A peer (synonymous with "client") is one instance of a BitTorrent client that does not have the complete file, but only parts of it.

Seeder: A seeder is a peer that has a complete copy of the torrent and still offers it for upload.

Torrent File: The torrent file contains metadata about all the files it makes downloadable, including their names and sizes and checksums of all pieces in the torrent.

Swarm: All peers (including seeders) sharing a torrent are called a swarm.

Tracker: A tracker is a server that keeps track of which seeds and peers are in the swarm. Clients report information to the tracker periodically and in exchange receive information about other clients to which they can connect.

The Prototype Components

The Client

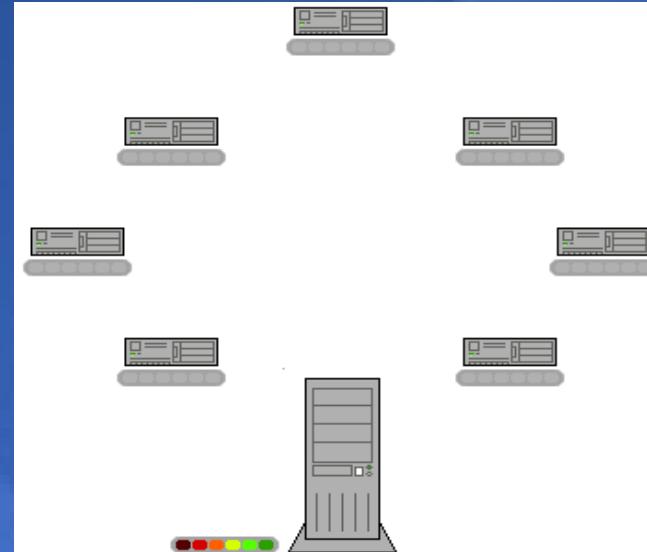
The Seeder *

The Tracker *

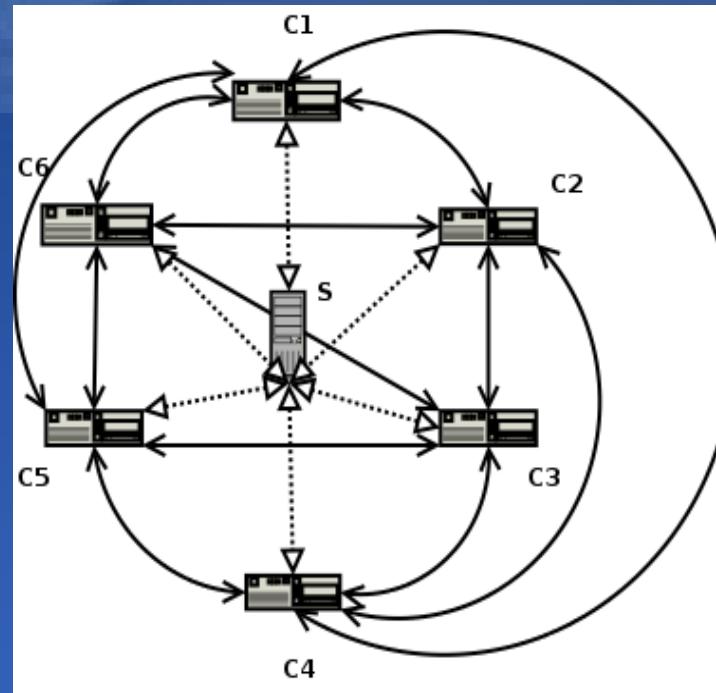
* Run as Stand-Alone Application (resident memory app)

The Client Application functionalities

- ◆ Databank Downloading using BitTorrent Protocol
- ◆ Torrent File generation for a local Bank repository
- ◆ Torrent File Publishing on the Tracker

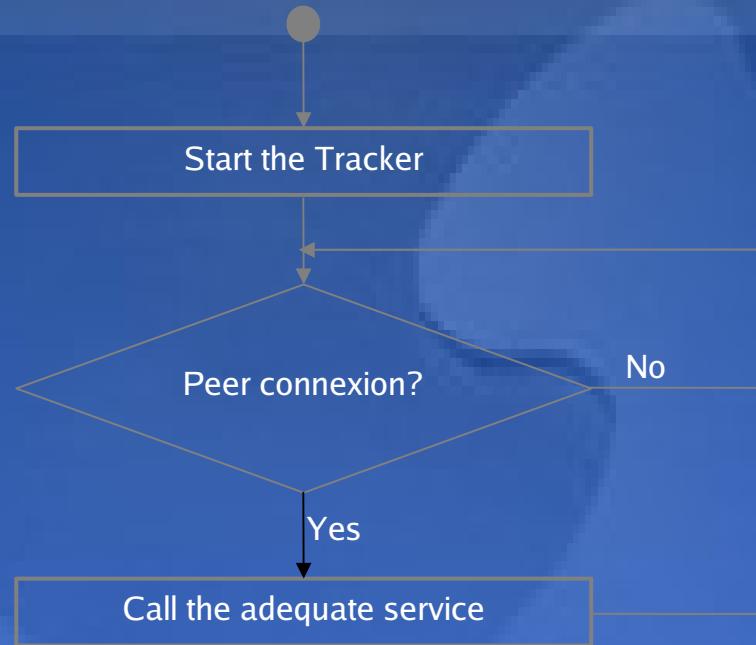


- ◆ DataBase Sharing



The Tracker App functionalities

- ◆ TrackerService for the sharing management
- ◆ TorrentServices for Torrent File Management

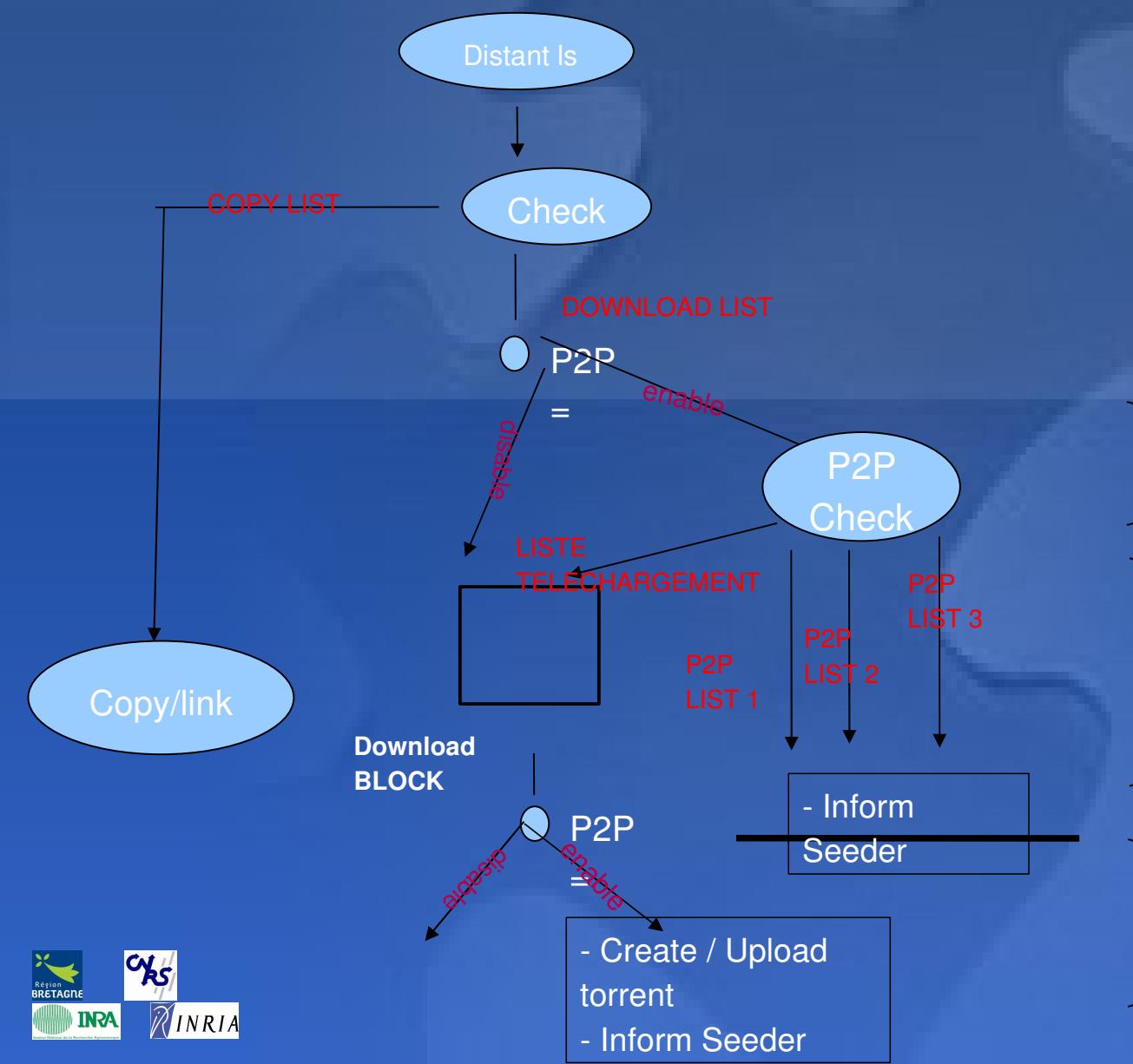


.Torrent File Content (MetaData)

```
d8:announce31:http://biomaj.genouest.org:80827:comment30:Propertiesfilegenbank
Release10:createdby9:HiKarkour4:info6:lengthi3583e4:name25:genbankRelease.properties12:pieceLengthi65536e6:pieces20:\D}ó~S~Eð÷~Am~@^F^Y^H~[~Hx@-Ñee
```

- ◆ Tracker URL
- ◆ File(s): Name, Length, Piece Length used, SHA-1 hash code for each piece
- ◆ Comments :
 - Server Name
 - Remote Dir
 - Remote Files & Remote Files Excluded
 - Attribute File List (size,date)
 - BioMAJ MetaData

The new Sync Cycle



From the Download List:

- New Download List
 - P2P / Torrent Lists

- N Workflows torrent (existing torrent)

- 1 actual Workflow
(Creating a new torrent)

- Inserting Torrent comments :
- Download LIST with file attributes
(date,size)
- SERVER (property)
- REMOTE.DIR (property) 14

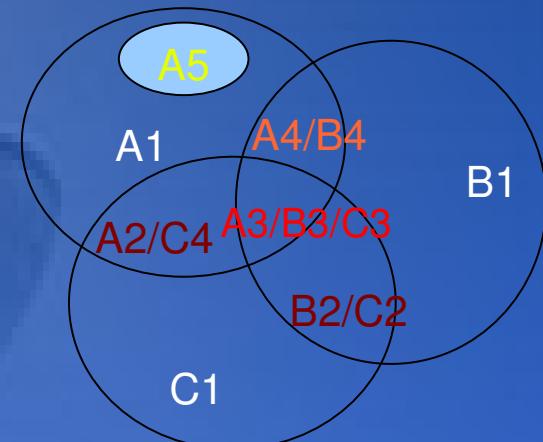
USE CASES

Case	Bank	Local Repository	Tracker	Choice
1	A*B*C	A{1,4,5}	B, C	C
2		A{1,3,4,5}	B, C	C
3		A{1,3,5}	B, C	B & C
4		A{1,2,4,5}	B, C	B C
5		A{1,4,5}	B	B & FTP(A2)
6		A{1,3}	B, C	B & C & FTP(A2)

Let 3 groups of files : A, B et C that have the following characteristics:

- Same ServerName&RemoteDir
- A, B and C have certain files in common

B and C are available via 2 torrents

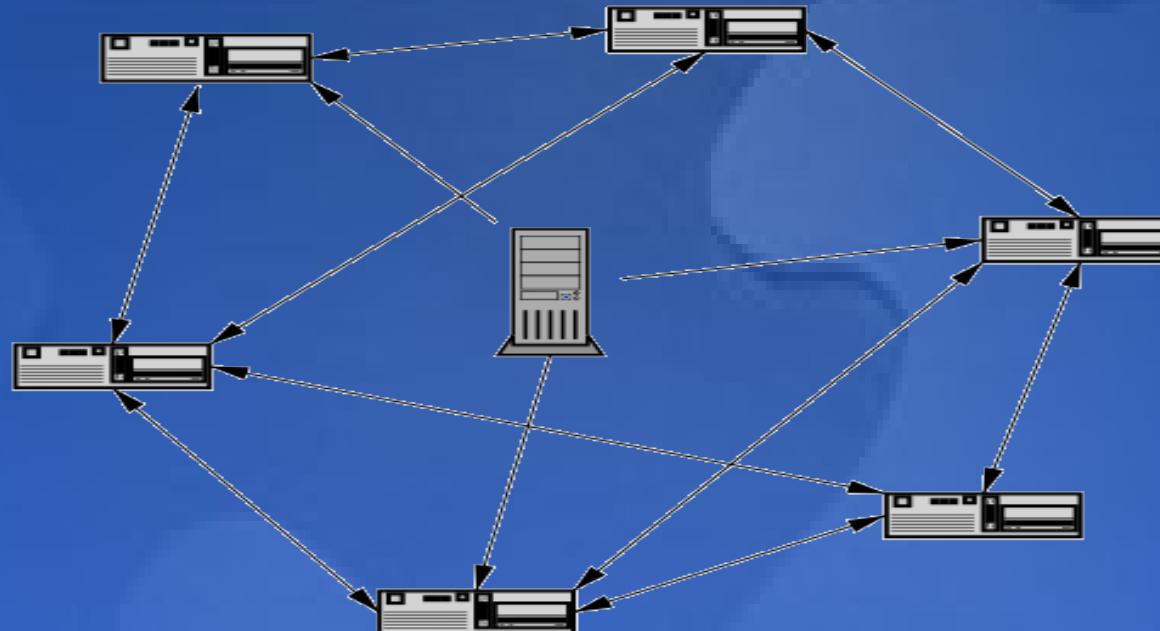


Security Features

- ◆ ACLS Access list at the Tracker Application Level

Notes :

- The Protocol tracker/client is based on HTTP messages
- The Protocol client/client is based on TCP



Roadmap

- ❖ Past :
 - Study of the available Peer-to-Peer protocols
 - Development of a Prototype
- ❖ Current :
 - Performance test on genbank (326GB) using the current Prototype
- ❖ Dev to come :
 - Tracker Modification (file search, file replacement)
 - Implementing the new Prototype to the BioMAJ WorkFlow
- ❖ Beta Version Release :
 - ~ March 2008
- ❖ Future work (Proposition) :
 - A “Cloud Computing” version running in a parallel grid environment with optional features (distributed data indexing, etc...)

Acknowledgements

- Sarah Dagher : Internship from ESIB-USJ
- DataMan & DocMen !

Q&A

<http://biomaj.genouest.org>

If what has been occurring in IT during the past decade can be classified as the '**information age**', then going forward it's going to be viewed more as the '**connection age**'

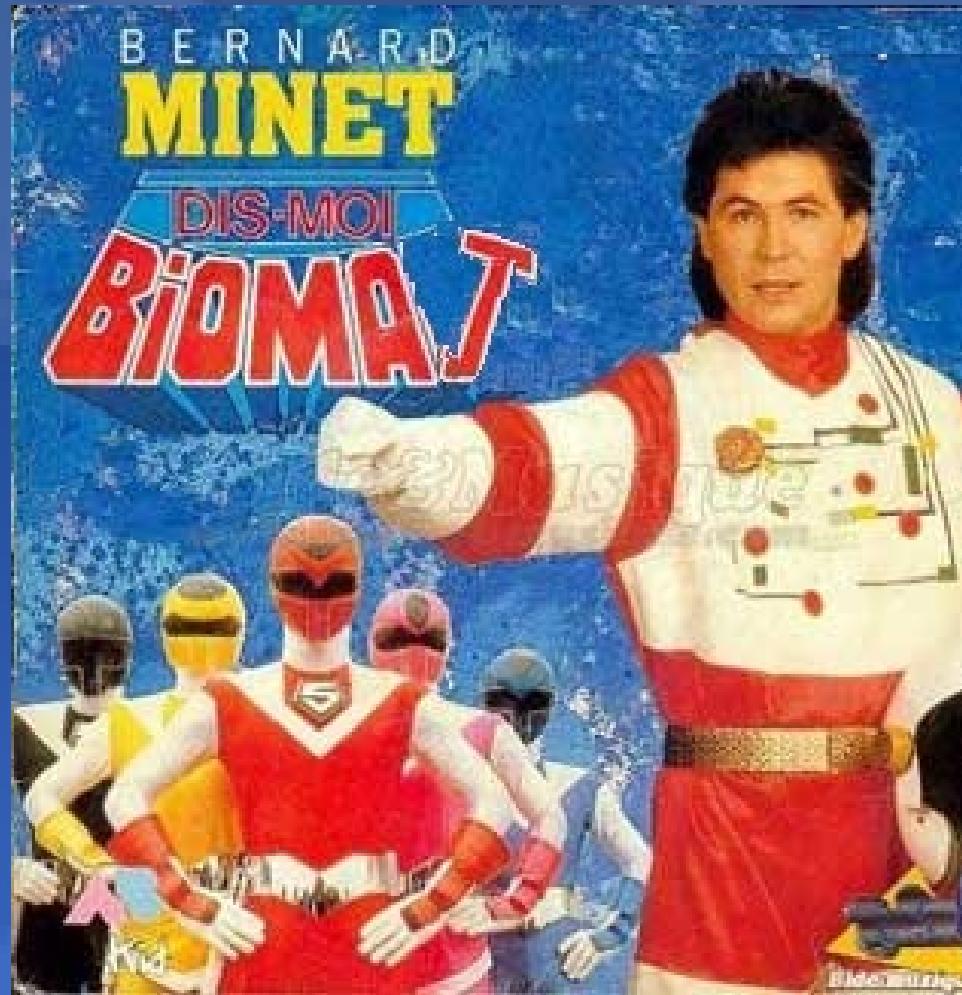
Ray Ozzie, CEO and chairman of Groove Networks (Year 2002)



Thank You All

P2P Works fine with Private DivX, Why shouldn't work with Public DataBanks!!!

BioMAJ



... The Return !