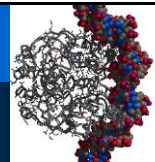# Exemple of
# Fast Translation Algorithm
## Object-oriented modeling and post-genomic
## biology : Programming Analogies

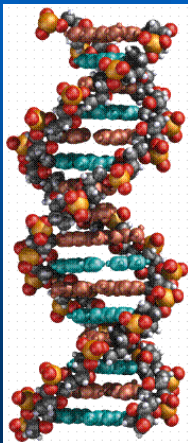**David Hill**
**Former Blaise Pascal University > UCA)**
**Université Clermont Auvergne**
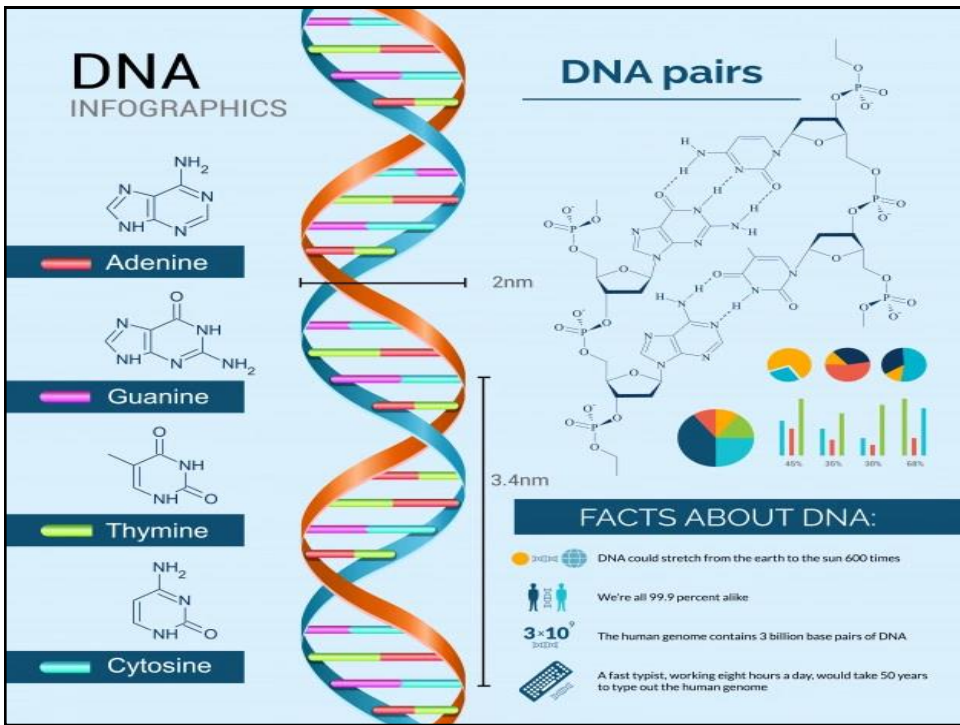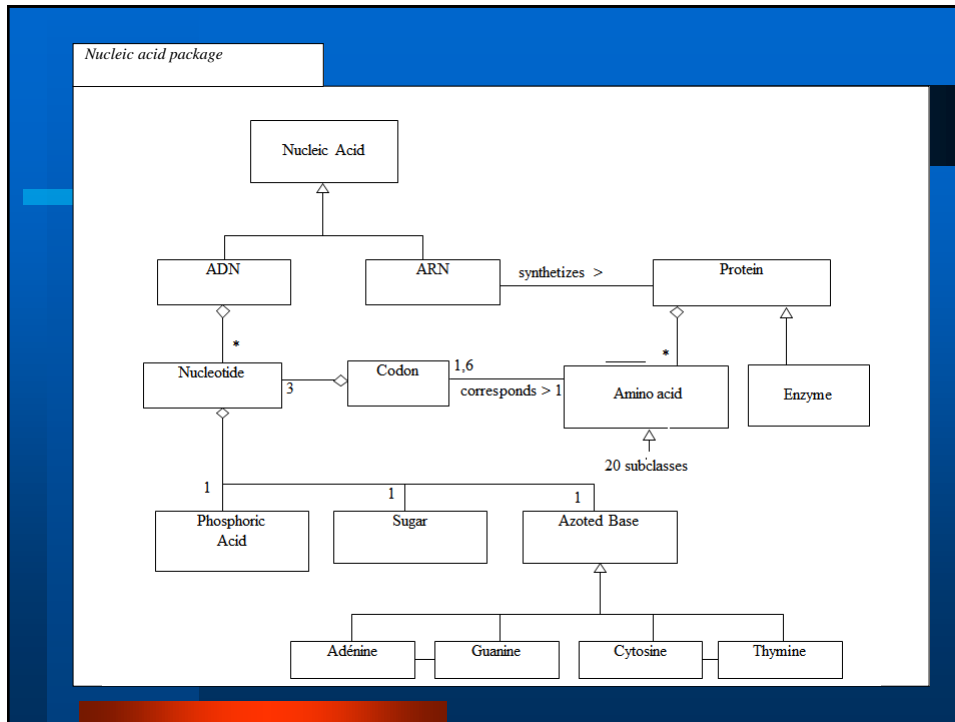**UMR CNRS 6158 : LIMOS**
**France**

# Key points

- **Molecular Biology basics**

- **Correspoding Classes in UML**

- **A translation algorithm**

- **Programming analogies**

- **Perspectives & Applications**
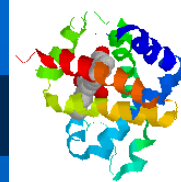
# Molecular Biology basics

- **DNA is made up of 4 individual chemical units called nucleic bases (A,C,G,T)**
- **DNA can act not only as a template for making copies of itself, but also as a blueprint for mRNA (messenger RiboNucleic Acid (repro. plan détaillé)**
- **The translation of mRNA into protein is the final major step in putting the information in the genome to work in a cell (in RNA Uracil U, replaces Thymine T)**
- **This kind of knowledge can be formalized by UML models and ontologies**
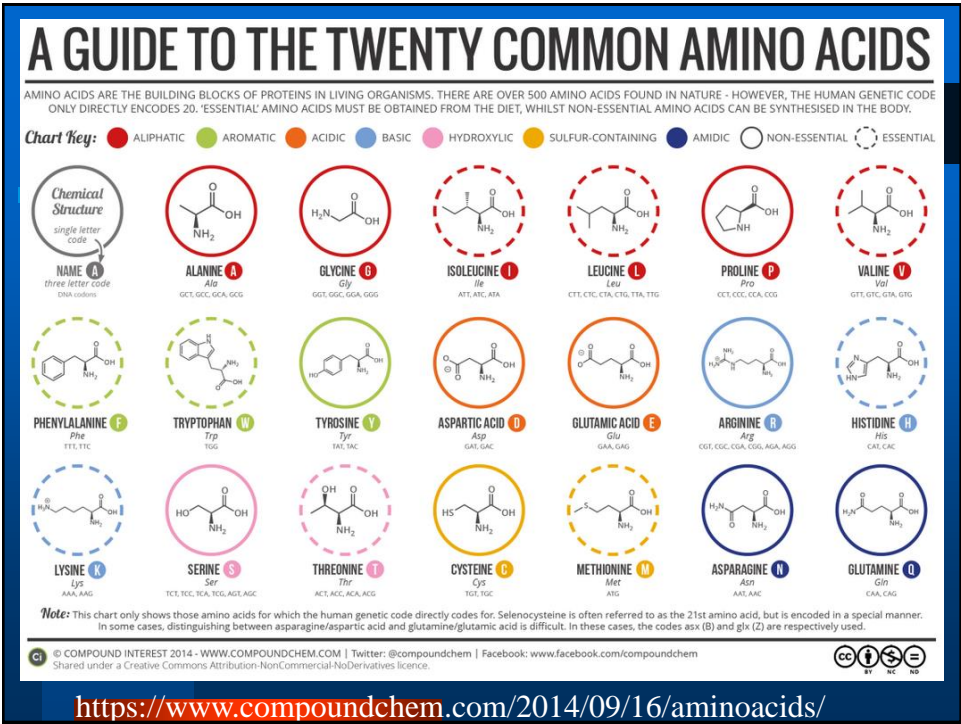
*Nucleic acid package*

Nucleic Acid

ADN — ARN — synthetizes > — Protein

Nucleotide — Codon — 1,6 — corresponds > 1 — Amino acid — Enzyme

3

* — 20 subclasses

1 — Phosphoric Acid

1 — Sugar

1 — Azoted Base

Adénine — Guanine — Cytosine — Thymine



# mRNA Translation process

- **A group of 3 bases is called a codon**
- **To each codon corresponds an amino acid**
- **There are 20 different amino acids**
- **We have coding redundancies since different codons give the same amino acid.**
  **64 codons => 20 amino acids**
- **The amino acid sequence specifies a protein (with start & stop sequences) or the enzymes.**

https://www.compoundchem.com/2014/09/16/aminoacids/

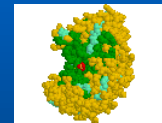# From Amino Acids to proteins

- **All the proteins that make up living organisms are huge molecules, but they're composed of tinier building blocks, known as amino acids.**
- **There are over 500 amino acids found in nature, yet, of these, the human genetic code only directly codes for 20.**
- **Every protein in your body is made up of some linked combination of these amino acids**
- **The previous graphic shows the structure of each, as well as giving a little information on the notation used to represent them.**

# Essentials or not ? (Covid again ?)

- **Broadly, these twenty amino acids can be sorted into two groups: essential and non-essential.**
- **Non-essential amino acids are those which the human body is capable of synthesising, whereas essential amino acids must be obtained from the diet.**
- **Some can also be termed 'conditionally essential', meaning that they may be needed from the diet during illness or as a result of health problems. This sub-category includes arginine, glycine, cysteine, tyrosine, proline, and glutamine.**
- **The essentials amino acids are histidine, isoleucine, leucine, lysine, methionine, phenylalanine, threonine, tryptophan and valine.**

# Start & stop codons…

- **Usual START codon : ATG (AUG)**

- **3 STOP codons : TAG, TGA, TAA**

- **Sequences between 2 stop codons are called open reading frames (ORFs), they are potentially coding for a protein**

- **With codons of 3 bases there are potentially 6 reading frames (3 on each complementary strands**
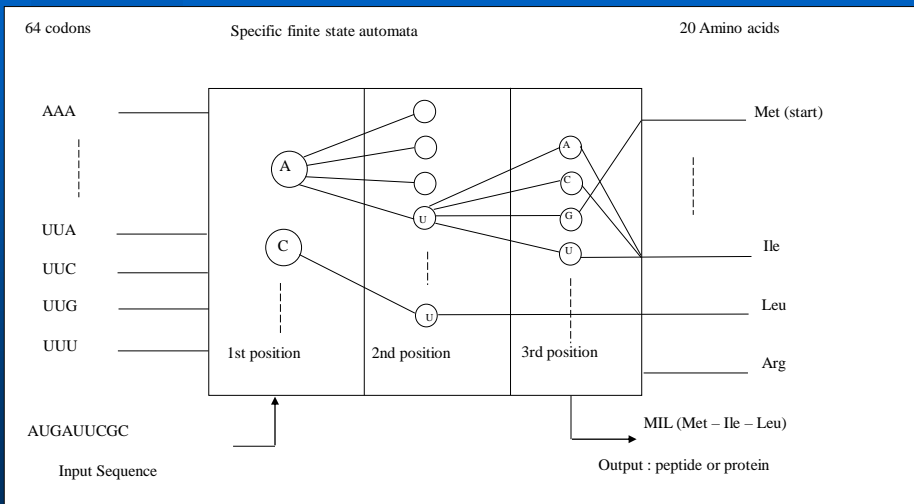
# The « genetic code »…

| 1st | 2nd | | | | 3rd |
|---|---|---|---|---|---|
| | **T** | **C** | **A** | **G** | |
| **T** | F Phe | S Ser | Y Tyr | C Cys | **T** |
| | F Phe | S Ser | Y Tyr | C Cys | **C** |
| | L Leu | S Ser | **Ter** | **Ter** | **A** |
| | **L Leu** | S Ser | **Ter** | W Trp | **G** |
| **C** | L Leu | P Pro | H His | R Arg | **T** |
| | L Leu | P Pro | H His | R Arg | **C** |
| | L Leu | P Pro | Q Gln | R Arg | **A** |
| | **L Leu** | P Pro | Q Gln | R Arg | **G** |
| **A** | I Ile | T Thr | N Asn | S Ser | **T** |
| | I Ile | T Thr | N Asn | S Ser | **C** |
| | I Ile | T Thr | K Lys | R Arg | **A** |
| | **M Met** | T Thr | K Lys | R Arg | **G** |
| **G** | V Val | A Ala | D Asp | G Gly | **T** |
| | V Val | A Ala | D Asp | G Gly | **C** |
| | V Val | A Ala | E Glu | G Gly | **A** |
| | V Val | A Ala | E Glu | G Gly | **G** |

# Translation algorithms

- **Bioinformaticians intensively use Translation of mRNA sequences and it is a very simple process compared to the sequence alignment problems (were the famous BLAST and FASTA programs are used)**
- **A fast translation process is modelled and implemented using :**
  - **Translation automata**
  - **(or even) Lockup tables**
- **Translation routines exist in every Bioinformatic package and library (GCG Translate, Bioperl translate,… )**

## A translation automata exploiting the genetic coding redundancies

64 codons | Specific finite state automata | 20 Amino acids

AAA

A

C

UUA

UUC

UUG

UUU

AUGAUUCGC

Input Sequence

1st position | 2nd position | 3rd position

A
C
G
U

Met (start)

Ile

Leu

Arg

MIL (Met – Ile – Leu)

Output : peptide or protein

## A very fast translation model

mRNA

Fasta or .nsq
File format

1 byte with
2 bits unused

Lookup in the binary
associative table if needed

3rd > 'G'    '10'

A        T

| | 1st base 00 | 2nd base 11 | 3rd base 10 |

$2^6 = 64$ combinations
6 bits codons

Simple computation :
1st base
<< 2 bits left shift
+ 2st base
<< 2 bits left shift
+ 3rd base

Lookup table [6 bits index]

Corresponding
Amino Acid : metionine

## Main translation code (C or C++)

```
// The following code supposes the initialization
// lookup tables. Thus majority of the code is composed of
// array initialisation + a simple computation of the
// binary index (given below for one codon)

index = Binary_Lkup   [ Bases[ i++ ] ];      // 1st base
index <<= 2;
index += Binaire_Lkup [ Bases[ i++ ] ];      // 2nd base
index <<= 2;
index += Binairy_Lkup [ Bases[ i++ ] ];      // Last base
                                             // for codon

// Look up of the corresponding Amino Acid
result = AminoAcid[index]
```

## Some translation times

| Linux Athlon 1.5 Ghz | CODING APPROACHES | | | | | |
|---|---|---|---|---|---|---|
| File size | Library routine | Better coding | Lookup tables | | Fastest algorithm | |
| # of input bases | | | Without opt. | With optimzation | Without opt. | With opt. |
| 1 016 684 | 8,0 | 2,8 | 1,1 | 0,52 | 0,1 | 0,09 |
| 2 033 351 | 16 | 5,4 | 2,0 | 1,0 | 0,2 | 0,2 |
| 4 066 684 | 32 | 11 | 4,1 | 2 | 0,4 | 0,37 |
| 8 133 351 | 63 | 22 | 8,3 | 4,1 | 1 | 0,6 |
| 16 266 684 | 126 | 43 | 17 | 8 | 1,9 | 1,2 |
| 32 533 351 | 251 | 86 | 33 | 16 | 3,2 | 2,5 |
| 64 066 684 | | | 66 | 32 | 6,6 | 5 |
| 130 133 351 | | | 132 | 65 | 12,6 | 9,5 |
| 260 266 684 | | | | | 25 | 19 |

**Translation times in seconds in 2002 !**

## Test at genome scale…
## on a basic PC and 2 Giga bases

- **On simple PC, the whole translation takes 3mins & 15 seconds (Including the reading of a 2Gb file (fasta format) and output of 670 Mb file ( in 2002! )**
- **The translation in itself takes 14 seconds for 2 Giga bases.**
- **Input / outputs now take 94% of the global response time**
- **=> if possible optimize input/ouputs with the Unix mmap function and specific HW (see code).**

## Program models & representations

*Finite automata*

*Fastest implementation*

100 %
Seem
Executable

6 % exec ?

94 % seem
non executable
=>
precomputed
data and array
initalisations

## Analogy between genetic « data » and optimized computer programs ?

- **The draft of the human genome sequence announces that only 1.1% of the genome is spammed by exons (coding regions), 24% in introns (non coding region) and 75% being (unknown) intergenic DNA**
- **Advances showed that some introns area could be coding in specific tissues**
- **What about the role of intergenic DNA, could it be used as data by the rest of the code ?**

## Applications @ LIMOS

- **Breast Cancer Research**
- **Design of oligonucleotides for DNA chips and microarrays**
- ***Study of the Encephalitozoon cuniculi parasite (smallest Eucaryote genome, sequenced at Blaise Pascal University)***
- **Bioremediation : design biological depollutant with Micro-arrays**

Research platform for DNA chips & microarrays.

probes (Human) (& mouse) → Preparation & amplification (PCR) → Spotting of DNA microarrays

Interpretation ← Analysis ← Hybridation ← mRNA



Object-oriented design of the data base (excerpt of UML class diagram)

# Example of input screens



# Verification of values in 96 holes plates

# DNA Microarray Image Analysis

- **Study of existing software**
  - **Scanalyze**
  - **Arrayvision**
  - **Genepix**
  - **Dapple**
  - **Jaguar**
  - **…**
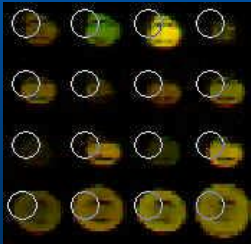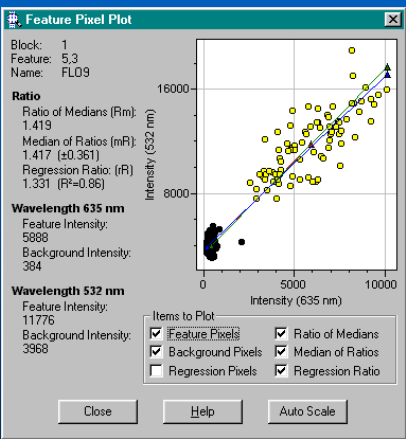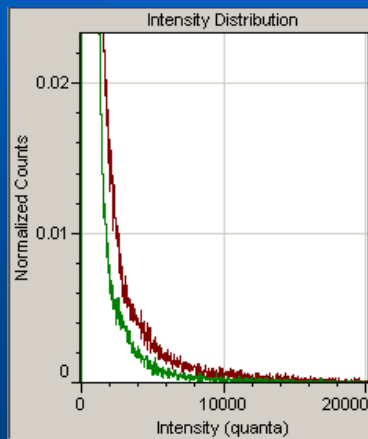- **Specific developements**



# Automatic Alignment

**Two algorithms are used :**
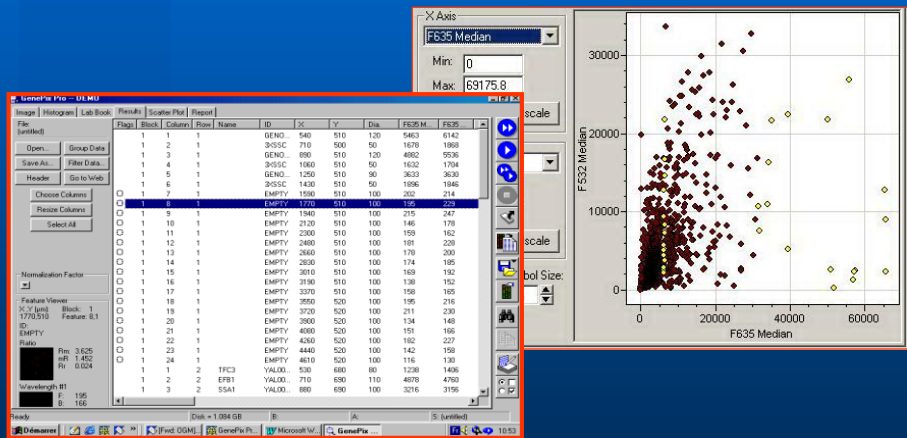
**The first giving a rough alignment**

**A second which will refine the first results**

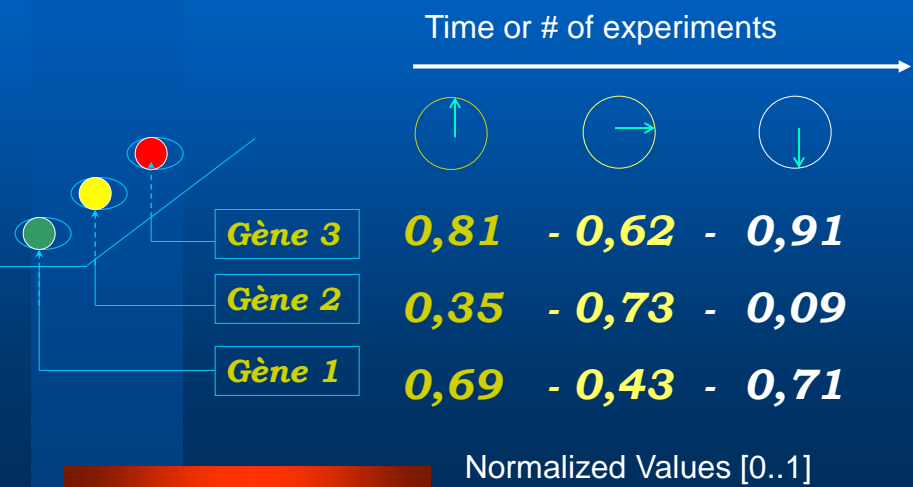Study of pixels distribution



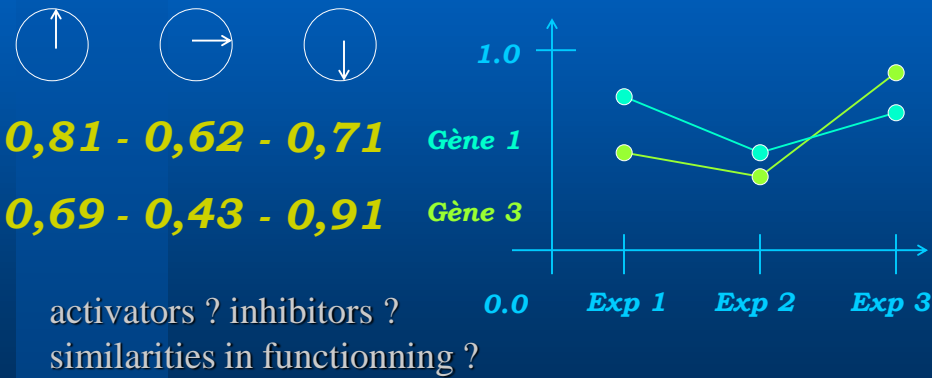Assignment of results enabling the caracterization of genes expression

# Classification & data-mining

- **Supervised & non supervised techniques :**
  - **Hierarchical clustering,**
  - **K-Mean and variants,**
  - **Aggregation algorithms**
  - **Self-Organizing Maps : SOM**
  - **Decision Trees (C5 &Classification with Regression Trees :CART)**
  - **Neural networks**
  - **Genetic algorithms**
  - **Information theory**
  - **Fuzzy logic,**
  - **Support Vector Machines (SVM)**
- **Data base techniques**
  - **Association rules (Apriori algorithm & its optimizations, LIMOS)**
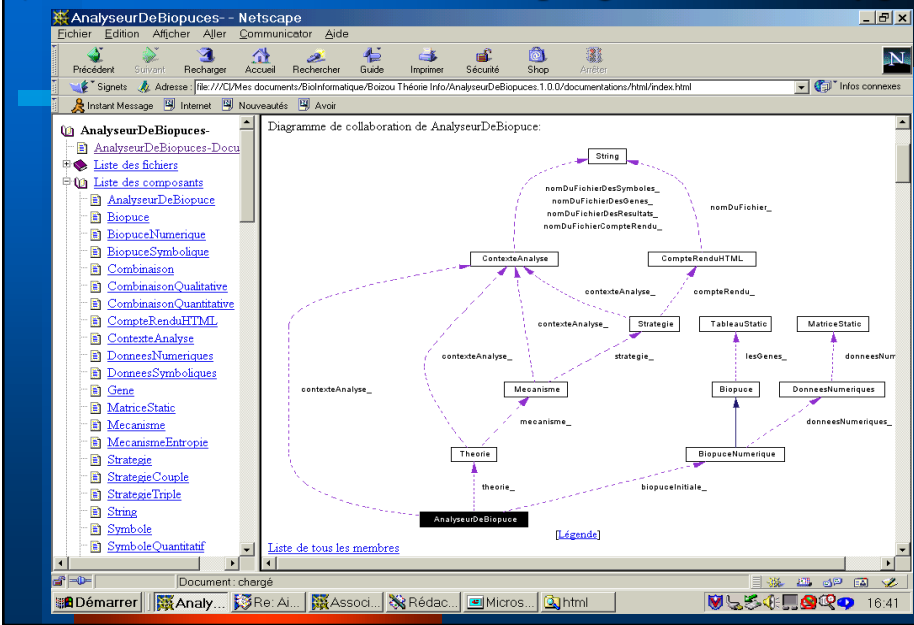  - **Functional dependencies  (DepMiner, LIMOS)**

# Resulting matrix to explore

Time or # of experiments

| | ↑ | → | ↓ |
|---|---|---|---|
| *Gène 3* | **0,81** | **- 0,62** | **- 0,91** |
| *Gène 2* | **0,35** | **- 0,73** | **- 0,09** |
| *Gène 1* | **0,69** | **- 0,43** | **- 0,71** |

Normalized Values [0..1]

## Research of similar or opposite behaviours



**0,81 - 0,62 - 0,71**   *Gène 1*

**0,69 - 0,43 - 0,91**   *Gène 3*

activators ? inhibitors ?
similarities in functionning ?
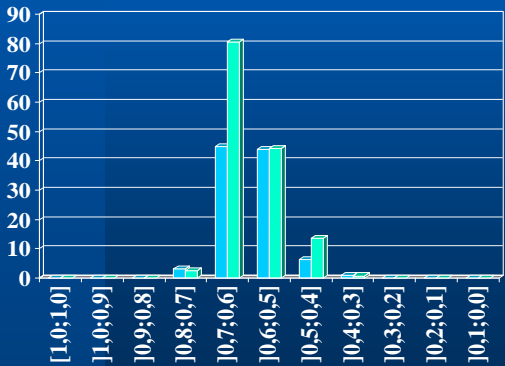
1.0

0.0   Exp 1   Exp 2   Exp 3

## OO Modeling and development
## (UML Documentation of C++ program with Doxygen)

## Study of the impact of noise in the observed results ...

### Sensibility to the number of experiments and to the number of classes

Example for :

100 random genes

10 Classes
20 Experiments

☐ Moyenne
☐ Variance



## Main partners

- **Computer Science & Modelling Laboratory (L.I.M.O.S)
  at Blaise Pascal University
  UMR CNRS 6158 – Now @ UCA & Clermont INP**

- **L.O.M. : Laboratoire d'Oncologie Moléculaire du
  Centre Jean-Perrin (Clermont-Ferrand)**

- **Cellular & Molecular Nutrition Unit at INRA (National
  Agronomy Research Institute)**

- **Old Protists Biology Laboratory (LBP) at Blaise Pascal
  University UMR CNRS**

## Perspectives

- **Study of SNP surprises ! (Single Nucleotide Polymorphism) and developement of simulations**
- **Study of the bioremediation process thanks to reverse transcription**
- **Development of gene prediction models at genome scale, based on the fast algorithm presented (discovering ORFs)**
- **Evaluation of classification techniques**

## From Web developer… Salon de la DATA …to Bioinformatics

https://www.youtube.com/watch?v=ocvSY74narw&ab_channel=SalonData

Résultats pour **salon de la data code source du vaccin** *Pfizer*   Essayez avec **salon de la data code source du vaccin Pfize**

A la découverte du code source du vaccin Pfizer/BioNTech
134 k vues • il y a 1 an

Salon Data

nouvelle technologie des **vaccins** ARNm est composé d'un **code source** d'un peu plus de 4000 caractères. nous ...

A la découverte du code source du vaccin Pfizer/BioNTech

Colin Cleary     @CColinCleary

Lire (k)

0:03 / 23:09

## Another possible optimization technique

- **Among the programming optimization techniques the unrolling technique is sometimes used**
- **Meta-programming can exploit this directly within the C++ language**
- **It helps in the coding of programs that will be unrolled at compilation time.**

## Unrolling with a C++ metaprogram

```cpp
template<bool> void IF(int tab[], int n)     { }
inline void IF<true>(int tab[], int n)       { tab[n] *= -1; }
inline void IF<false>(int tab[], int n)      { tab[n] = 0; }

template<int N>
inline void FOR(int tab[])
{
    IF<N%2>(tab,N);
    FOR(N-1)(tab);
}
inline void FOR<-1>(int tab[])               { }

const int N = 10;
void main()
{
    int * tab, i ;

    tab = new int[N];
    for(i = 0 ; i < N ; i++) cin >> tab[i];
    FOR<N-1>(tab);
    for(i = 0 ; i < N ; i++) cout << tab[i];
}
```

## THANKS FOR YOUR ATTENTION

**More to come**
**Invitation to C++ TMP ZZ3**