An INFOrmational GENOMICS approach (based on joint work with V. Bonnici, A. Castellini, V. Manca)

Dr Giuditta Franco

Department of Computer Science, University of Verona, Italy

Dr Giuditta Franco An INFOrmational GENOMICS approach

イロト イポト イヨト イヨト

Computational genomics

Study of genomes is referred to as *genomics* - while *genetics* concerning the study of (single of groups of) genes.

Computational/statistical genomics include gwas (genome wide association studies), wgs (whole genome sequencing algorithms), *advanced data structures (suffix trees, arrays, BWT)* and big data approaches (to process massive amount of data).

Local sequence (similarity) analysis is traditionally performed by **alignment-based methods** (FASTA, BLAST). Traditional **alignment-based methods** are applied for multiple sequence analysis and comparison, with phylogenetic and pharmacogenomics applications.

・ロット (雪) (日) (日)

Alignment free methods

Recent (ten years) **alignment-free methods**¹ in computational genomics focus on a systemic view, based on empirical studies of frequencies of DNA *k*-mers (genome factors of length *k*), out of whole genomes/proteomes ², where FL and information theories are applied.

For a fixed length code, empirical frequencies are computed on *k*-mers as normalized multiplicities.

²set of proteins possibly expressed by a cell, tissue, organisma > < a > a < o < c

¹Vinga et al. '03, Fofanov et al., 04, Chor et al. '09, Searls '10

Dictionary-based genomics

As a *very long* string over an alphabet of four/five symbols, a genome may be seen as a text book, which language has to be completely deciphred. Information is comprised in words, *lineraly* arranged by unknown synthax and semantics.

Sequences become **bags of words**, used in computer vision, for document classification by clustering of strings which keep only their multiplicity/feature (*no order*, no grammar).

Investigation and comparison (e.g., similarity analysis) by **dictionaries** or multisets. Models common to natural language processing, information retrieval (IR), computational linguistics. This abstraction excludes tandem repeats.

< ロ > < 同 > < 三 > < 三 > 、

Informational genomics: main ingredients

Genomes are information source, whose words are random variable values, data encodings, and a probability distribution is given by frequencies ³.

Infogenomics employs methods from FLT ⁴ and IT to define genomic indexes, dictionaries, distributions, elongation and segmentation methods, sequence similarity measures, gene networks ⁵.

- ⁴V. Brendel at al. Nucleic Acids Research 12, 1984.
- ⁵A. Castellini et al. BMC Genomics 2012, Nat Comp 2015.

³Sims et al. PNAS, 2008. Robins et al. Journal of Bacteriology 2005

Lecture outline

An informational dictionary-based alignment-free approach:

- Genomic dictionaries, distributions, indexes.
- Multiplicity-comultiplicity genomic profiles, and UCE (ultraconserved elements) by dictionary intersections
- Systematic analysis of repeats variation in number, length, multiplicity, and localization (inside, outside genes)
- Genomic Recurrent Distance Distribution (**RDD**), and cluster analysis on repeat sharing **gene networks**.

Software Tools

< ロ > < 同 > < 三 > < 三 > 、

Our work focused on real genomes

Dictionaries of 60 specific genomes, deeper analysis on 12



Organism Genome	Length	Genes	Туре
Nanoarchaeum equitans	490,885	585	Minimal archaeum
Mycoplasma genitalium	580,076	476	Minimal bacterium
Mycoplasma mycoides	1,211,703	1,016	Venter's experiment bacterium
Haemophilus influenzae	1,830,138	1,717	First sequenced bacterium
*Escherichia coli	4,639,675	4,685	Bacterium model (K-12)
*Pseudomonas aeruginosa	6,264,404	5,566	Ubiquitous bacterium
*Saccharomyces cerevisiae	12,070,898	6,275	Unicellular eukaryote (Yeast)
Sorangium cellulosum	13,033,779	9,700	Longest genome bacterium
Homo sapiens chr. 19	63,800,000	2,066	Highest gene density H. chr
*Caenorhabditis elegans	100,267,632	19,000	Worm (around 1000 cells)
*Drosophila melanogaster	129,663,327	14,000	Insect (fruit fly)
Homo sapiens chr. 1	247,000,000	3,511	Longest Human chr. 😑 🔊 🤉 🔿

Dr Giuditta Franco

An INFOrmational GENOMICS approach

Genomic distributions nformational indexes

Basic definitions

Alphabet $\Gamma = \{a, t, c, g\}$, genome $G \in \Gamma^*$.

• $D_k(G)$: genomic *k*-dictionary collecting all *k*-mers in *G*, $F_k(G)$ of forbidden *k*-words, which do not occur in G

Ex: *attaggatcttaat* has nine 2-words: six occurring once (aa, ag, tc, ct, ga, gg), two occurring twice (ta, tt), one (at) occurring 3 times, and seven 2-forbidden.

- *H_k(G)*: dictionary of hapax *k*-words, occurring once in G,
 R_k(G) the set of *k*-repeats, occurring in G at least twice
- *T_k(G)* is the *k*-factor multiset of G: a function over *D_k(G)* associating each string to its number of occurrences in G

イロン 不良 とくほう 不良 とうほ

Genomic distributions Informational indexes

Multiplicity-comultiplicity profiles

Multiplicity-comultiplicity 2-distribution diagram for the sequence *attaggatcttaat*



Genomic distributions Informational indexes

Genomic profiles



Dr Giuditta Franco An INFOrmational GENOMICS approach

э

Genomic distributions Informational indexes

A related string problem

Related to the *word assembly problem*: genome reconstruction method, from a dictionary (examers), or a mult-comult diagram.



M. genitalium, H. influenzae, E. coli, D. melanogaster.

These studies could improve existent genome reconstruction algorithms, by estimations of reads length repeatability.

Genomic distributions Informational indexes

Comparison among genomic profiles - examples





Above, *E. coli*, *H. influenzae*, *P. aeruginosa* are compared by their normalized genomic (6-)profiles.

Aside, *M. genitalium*, *E. coli*, *S. cerevisiae*, *H. sapiens chr 19*, are compared with one random permutation (in red) by the multiplicity-comultiplicity profiles.

< < >> < <</>

Genomic distributions Informational indexes

Fairly occurring motifs common to genomes (UCE)



Genomic distributions Informational indexes

Sizes of k-dictionaries: : real vs random genomes



E. coli 's genomic dictionaries *H. sapiens chr19* 's dictionaries



Dr Giuditta Franco

An INFOrmational GENOMICS approach

Genomic distributions Informational indexes

Sizes of k-dictionaries: : virus and bacterium

BeanGoldenYello's dicts







Burkholderia 's dicts



Dr Giuditta Franco An INFOrmational GENOMICS approach

Genomic distributions Informational indexes

A related open problem

Observed curves of $|D_k|$ (of $|H_k|$, $|R_k|$, and F_k) exhibit a similar shape for some genomes having a sensibly different length.

Open problem: the discovery and comprehension of some rule explaining the empirically evident relationship among genome length n, factor length k, and k-dictionaries cardinality.

イロト イポト イヨト イヨト

Genomic distributions Informational indexes

Phase transitions for hapax/repeat cardinality ratio

$$DT_k = rac{|D_k(G)|}{|T_k(G)|}$$
 (k-lexicality), $HR_k = rac{|H_k(G)|}{|R_k(G)|}$

Genomes	DT ₆	DT ₁₂	DT ₁₈	HR ₆	HR_{12}	HR ₁₈
N. equitans	0.008	0.87	0.99	1.468 ×10 ^{−3}	8.39	737.25
M. genitalium	0.007	0.85	0.98	8.65×10^{-3}	7.175	91.44
M. mycoides	0.003	0.53	0.81	9.661 ×10 ⁻³	2.169	12.33
H. influenzae	0.002	0.81	0.98	0	5.240	88.93
E. coli	0.0009	0.74	0.98	÷	3.331	115.84
P. aeruginosa	÷	0.47	0.98		1.564	93.76
S. cerevisiae		0.54	0.95		1.518	58.67
S. cellulosum		0.29	0.96		0.993	41.12
H. sapiens chr19		0.19	0.75		0.455	17.27
C. elegans		0.13	0.89		0.286	19.86
D. melanogaster		0.12	0.90		0.114	32.56

 $|T_k|$ = number of *k*-mers counted with their multiplicity (|G|-k+1). HR18 explains why microarrays work well,

Genomic distributions Informational indexes

Analysis of (relatively) long repeats reduction



Genomic distributions Informational indexes

Genomic Dictionaries

- D(G) = {G[i,j] |1 ≤ i ≤ j ≤ |G|} (square dim. w.r.t. |G|)
- D_k(G) = D(G) ∩ **Γ**^k
- L included in D(G) is a dictionary of G
- A position p of G is m-covered in D if there are m words G[i,j] of D with i ≤ p ≤ j (positional coverage)
- D covers G if every position of G is k-covered with k ≥ 1 by D (lexical coverage)
- D minimally covers G if D covers G and no D' included in D covers G
- G is D-segmentable if G belongs to D*

Slide courtesy of prof. Vincenzo Manca, Univ. Verona, IT

・ロット (雪) (日) (日)

Genomic distributions Informational indexes

Basic Genomic Indexes

-	LG	Logarithmic Length
-	LX_k	k-Lexical Multiplicity
-	MFL	Minimal Forbidden Length (MCL = MFL -1)
-	MRL	Maximum Repeat length (+1 = all-hapax lub least
		upper bound)
-	MHL	Minimum Hapax Length (-1 = all-repeat <i>glb</i> greatest
		lower bound)
-	COV	Coverage percentage (w.r.t. a dictionary)
-	POV	Positional coverage (w.r.t. a dictionary)
-	E _k (G)	Empirical k-Entropy
-	ED _k (G_1, G	_2) k-Entropic Divergence
-	IDG	Inverse de Brujin graph indexes

Slide courtesy of prof. Vincenzo Manca, Univ. Verona, IT

・ロン ・四 と ・ ヨ と ・ ヨ と …

= nac

Genomic distributions Informational indexes

Informational indexes

Minimal Hapax Length (genome itself is an hapax, any word including an hapax is an hapax), *Maximal Repeat Length* (any subword of repeat is a repeat), *Minimal Forbidden Length*

			[Ra]	
Genomes	MF	MR	40.000	
			70.000 Mycoplasma myrcides 36.000	Escherichia coli K-12
N. oquitana	6	120	0000	
N. equitaris	0	139	50.000 25.000	
M. genitalium	6	243	40,000 20,000	
M. mycoides	6	10.963	2000 10,000	
H influenzae	7	5 563	10,000 5,000	
	, <u>,</u>	0,000	0 0 2,000 4,000 6,000 0,000 10,000 0 0 000 1,000 1,000	2.000 2.500 k
E. COli	/	2,815	(Ra) (Ra)	
P. aeruginosa	8	5.304	200.000	
S coroviciao	o i	0 275	60.000 Previoenceas serviginesa 180.000	Saccharomyces cerevisiae
S. Cerevisiae	9	0,375	50.000 140.000	
S. cellulosum	7	2,720	40.000	
H. sapiens chr19	9	2,247	30.000 80.000	
C. elegans	10	38,987	20.000	
Dimelanamentar		00,000	10.000 20.000	
D. meianogaster	11	30,892	0 1,020 2,020 3,003 4,020 5,020 0 0 2,020 4,030	6.000 8.000
			- k	k

イロト イポト イヨト イヨト

Genomic distributions Informational indexes

MRL of virus and bacteria

97% viruses has MRL ranging 10-10³, 93% bacteria 100-10⁴



Slides courtesy of Dantoni-Mancini-Sartea, Univ. Verona, IT

Dr Giuditta Franco An INFOrmational GENOMICS approach

A B > A B >

Genomic distributions Informational indexes

Human Microbiome MRL



Genomic distributions Informational indexes

MHL and MFL for viruses

Major part of viruses has MHL and MFL between 4 and 6.



Slides courtesy of Dantoni-Mancini-Sartea, Univ. Verona, IT

Dr Giuditta Franco An INFOrmational GENOMICS approach

イロト イポト イヨト イヨト

Genomic distributions Informational indexes

MHL and MFL for bacteria

Major part of bacteria has MHL and MFL between 6-8.



Slides courtesy of Dantoni-Mancini-Sartea, Univ. Verona, IT

Dr Giuditta Franco An INFOrmational GENOMICS approach

イロト イポト イヨト イヨト

Genomic distributions Informational indexes

Human Microbiome: MHL and MFL confirmed



Slides courtesy of Dantoni-Mancini-Sartea, Univ. Verona, IT

Dr Giuditta Franco An INFOrmational GENOMICS approach

・ロト ・ 同ト ・ ヨト ・ ヨト

э.

Genomic distributions Informational indexes

Informational indexes: some properties

For any genome G (of length *n*):

$$H_k = \Gamma^k \cap H, R_k = \Gamma^k \cap R \Rightarrow D_k = H_k \uplus R_k, \Gamma^k = D_k \uplus F_k$$

$$AR_k = \frac{|T_k \setminus H_k|}{|R_k|}$$
 average k-factors repeatability

 $S_n = \lfloor \lg_4 n \rfloor - MF + 1$ factor length selectivity

$MFL \leq MHL + 1$

・ロト ・ 同ト ・ ヨト ・ ヨト

Э

Repeat sharing gene networks RDD distributions Software Tools

Importance of genomic repeats

Repeats located in genes (insertion of ALU sequences) have a relevant role for survival and evolution of primates, as they may alter the rate of protein production

Correlation with number of mutations in (encoding regions of) many cancer-related genes (such as PTEN)

Tumors with microsatellite instability (MSI), characterized by massive instability in repeated sequences, display a strong microsatellite mutator phenotype (MMP)

Fragile X syndrome (a genetic disorder inducing mental retardation) is associated to the expansion of CGG affecting the (FMR1) gene on the X chromosome

・ロット (雪) (日) (日)

Introduction Repeats analysis Software Tools

Repeat localization (inside, outside, across genes)

Percentage of k-repeat occurrences, with the length k



Dr Giuditta Franco

Repeat sharing gene networks RDD distributions Software Tools

Repeat sharing gene networks

Def. A *k*-parametrized, labeled graph $G_k = (V_k, E_k)$, where:

- V_k are nodes associated to the genes of the organism,
- two genes are connected (in E_k) if they share at least one k-repeat. The set of shared k-repeats is the edge label.

By sketching the *k* value, nodes and edges in G_k decrease (as isolated nodes are deleted), until the network disappears⁶. A break-point (*k* = 18) where the networks pass from having a big connected component to being a set of "vanishing" clusters.

⁶In *N. equitans*, G_{54} is empty, while MR = 139

Repeat sharing gene networks RDD distributions Software Tools

Genetic repeats within a specific genome

N. equitans 's gene networks G₁₄ - G₁₈



Repeat sharing gene networks RDD distributions Software Tools

Another example: *E. coli* 's gene network G_{18}



Dr Giuditta Franco An INFOrmational GENOMICS approach

Investigation on repeat sharing gene networks

Max degree, max labels weight, number of cliques (its variation with k) - results on *N. equitans, E. coli, S. cerevisiae*

Highly connected genes, for k long enough, have similar functionality and turned out involved in important biological pathways, such as DNA repair and replication.⁷

Gene length compared with repeat length k, to measure edge significance (genes shorter than 100 encode for tRNA)



⁷A. Castellini et al., Natural Computing 2015

An INFOrmational GENOMICS approach

A B + A B +
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

Repeat sharing gene networks RDD distributions Software Tools

Basic E.coli network analysis





・ロト ・ ア・ ・ ヨト ・ ヨト

Э

Dr Giuditta Franco An INFOrmational GENOMICS approach

Repeat sharing gene networks RDD distributions Software Tools

Similar network patterns

Network curves observed for *N. equitans, E. coli, S. cerevisiae* do not depend on genome length or gene coverage.



Dr Giuditta Franco An INFOrmational GENOMICS approach

イロト イポト イヨト イヨト

Introduction Repeat sharing gene networks Genomic dictionaries RDD distributions Repeats analysis Software Tools

Complete clusters (N. equitans, E. coli, S. cerevisiae)

For longer repeats, gene networks tend to aggregate in cliques.



Repeat sharing gene networks RDD distributions Software Tools

Examples of cliques (in E. coli)



Dr Giuditta Franco An INFOrmational GENOMICS approach

・ロト ・ 同ト ・ ヨト ・ ヨト

ъ

Repeat sharing gene networks RDD distributions Software Tools

Repeat clique analysis

According to the range values of k, we may deduce:

 $k = 1, \ldots, 12$: repeats are completely random

k = 13, ..., 20: some repeats are present only in couples of genes, only few have a biological role (14-15: first repeats present even in non-tRNA genes)

k > 21: Repeats (all, for k > 40) have a biological role, they belong to paralogous genes and have a same reading frame for protein translation.

イロト イポト イヨト イヨト

Repeat sharing gene networks RDD distributions Software Tools

Clique analysis - a few observations

 Relatively few genes are involved in cliques, and the number of cliques varies similarly in the three organisms.

In every cliques, there is at least one repeat common to all



◇ Such a repeat encodes for a protein/enzyme core, and has the same reading frame in *all* the genes where it occurs.

Almost all cliques are composed by paralogous genes.

Repeat sharing gene networks RDD distributions Software Tools

Minimal k-RDD

Let $\alpha \in D_k(G)$ such that:

 $--\alpha$ $--d_1$ $--\alpha$ $--d_2$ $--\alpha$ $--d_1$ $--\alpha$ $--d_3$ $--\alpha$ ---

where no α occurs between two consecutive α

K-RDD(α , G) = d₁ -> n₁, d₂ -> n₂, d₃ -> n₃, ---

If G is random, for a "not too short, but long enough", K-RDD(α , G) is geometric/exponential, according to probability theory.

Slide courtesy of prof. Vincenzo Manca, Univ. Verona, IT

◆□▶ ◆□▶ ◆三▶ ◆三▶ ・三 ・ のへで

Repeat sharing gene networks RDD distributions Software Tools

ATG in human chr 22 vs extimated Exp Distr



Dr Giuditta Franco An INFOrmational GENOMICS approach

Repeat sharing gene networks RDD distributions Software Tools

ATG in human chr 22 and sequences at distance 81





Dr Giuditta Franco An INFOrmational GENOMICS approach

Repeat sharing gene networks RDD distributions Software Tools

ATC in E. coli: peaks with non-repetitive elements





< 🗇 ▶

▶ < ∃ >

Repeat sharing gene networks RDD distributions Software Tools

Average RDD for k=4, Emiliana huxley virus86



Dr Giuditta Franco An INFOrmational GENOMICS approach

▶ < Ξ >

Repeat sharing gene networks RDD distributions Software Tools

RDD in exonic regions, k=3

Peaks at distance 3 disappear in transcripts (introns + exons) - they are a *codonic language*



Repeat sharing gene networks RDD distributions Software Tools

Repetitive elements in ex. regions distance multiple 84



k=6



Dr Giuditta Franco An INFOrmational GENOMICS approach

Repeat sharing gene networks RDD distributions Software Tools

Importance of grouped occurrences

Recurrence is a peculiar feature of words

Occurrences of words semantically relevant for a text tend to be grouped around some points

 $RDD(\alpha)$ is a measure to visualize the non-randomicity of the word α distribution

・ロット (雪) (日) (日)

Repeat sharing gene networks RDD distributions Software Tools

Segment multiplicity and recurrent words

A **Bernoullian genome** is generated by means of casual extraction (with insertion after extraction) from an urn containing balls of four colours.

Segment multiplicity. Let us consider the genome as the concatenation of equal length segments. Given a word α , this distribution assigns to each *n* the number of segments where α occurs *n* times.

・ロット (雪) () () () ()

Repeat sharing gene networks RDD distributions Software Tools

Bernoullian (random) genomes

In a Bernoullian genome, such distribution (normalized, with the total num of segments) of word frequencies follows a Poisson prob distribution (of a certain variance λ): $Pr(X = k) = \frac{\lambda^k e^{-\lambda}}{k!}$



< ロ > < 同 > < 三 > < 三 > 、

The **waiting time** follows the Poisson law: the distance between two consecutive occurrences of a given α is an exponential of parameter *h*, $f(x) = he^{-hx}$, $x \ge 0$, for some *h*.

Repeat sharing gene networks RDD distributions Software Tools

Tandem repeats investigation

RDD (recurrence distance distribution): given a sequence α , to each *n* it is assigned the number of times that α occurs at distance *n* from its previous occurrence. Once normalized, the distribution above may be compared with (as corresponds to) the **waiting time** associated to a Poisson process.

If x is the distance between two occurrences of a given string, P(x) = f(x, h) is the number of times two occurrences appear at distance x. The probability of occurring at distance x is the ratio between the number of times it recurs at distance x and the multiplicity of α in G. A curve average may be then computed, over all sequences having the same length of α .

・ロト ・ 同ト ・ ヨト ・ ヨト

Repeat sharing gene networks RDD distributions Software Tools

Information Correlation and RDD in Genomes

- Trifonof et al. : DNA correlation periodicities, 1980
- Shepherd : DNA periodicities in coding regions, 1981
- Eigen et al. : periodicity in Transfer-RNA, 1981
- Fickett :1982 non min. RDD periodicity in coding regions, 1982
- Li : Mutual information in DNA Strings, 1990
- Herzel et al. : Measuring DNA correlations, 1990
- Li internal correlation in DNA, 1997
- Herzel-Weiss-Trifonof : 10-11 Periodicity, 1999
- Afreixo : 1-RDD min. 2009
- Bastos : 2-RDD min. 2011
- Carpena et al. RDD in keywords finding (non DNA), 2009-20013
- Computational Chemistry 2014

Slide courtesy of prof. Vincenzo Manca, Univ. Verona, IT

・ロット (雪) (日) (日)

E DQC

_

Repeat sharing gene networks RDD distributions Software Tools

Software for massive computations



Dr Giuditta Franco An INFOrmational GENOMICS approach

э

Repeat sharing gene networks RDD distributions Software Tools

Infogenomics Explorer

Visualization and exploration of informational indexes by means of a $Qlik^{(R)}$ View application called InfoGenomics.



Dr Giuditta Franco

An INFOrmational GENOMICS approach

Repeat sharing gene networks RDD distributions Software Tools

IGtools

Interactive graphical interfaces and CLI (batch analyses). Advanced data structures and algorithms, for real genomes.



Dr Giuditta Franco

An INFOrmational GENOMICS approach

Repeat sharing gene networks RDD distributions Software Tools

Conclusions

A method

- to represent and compare genomes (genomic profiles, dictionary intersections)
- to describe a genome by numerical information: statistics (amount, multiplicity, localization of repeats), informational index vectors
- to study gene networks, in general and along with a complete clusters analysis
- to find tandem repeats, and "good" genomic dictionaries.

・ロット (雪) (日) (日)

Repeat sharing gene networks RDD distributions Software Tools

Biological and computational annotations

Biological annotations of DNA elements, by the functional role they play in regulatory mechanisms (biological semantics).

Numerical annotation of genomic words, by the position in the genome, the total number of occurrences, the occurrences lying inside, outside, or between encoding regions, its CpG content, and more sophisticated informational indexes of text analysis.⁸

⁸G. Franco, *Perspectives in computational genomics*, Discrete and Topological Models in Molecular Biology, 2014.

Repeat sharing gene networks RDD distributions Software Tools

Open problems

♦ How information is structured within genomes? Finding a "good" genomic dictionary (cardinality, word length k, sequence and average positional "coverage")

 Regularity properties, informational indexes characterizing (classes of) genomes: e.g., MRL, MFL

 Genome discrimination methods, namely for phylogenetic or medical purposes.

・ロット (雪) (日) (日)

Repeat sharing gene networks RDD distributions Software Tools

Applications beyond genomics

Phylogenetic methods

Cyber virus code analysis (by breaking down malware signatures)

Computational linguistics

The dictionary parameter k has to be appropriately chosen for each model (a range of k may be possibly required)

< ロ > < 同 > < 三 > < 三 > 、

Repeat sharing gene networks RDD distributions Software Tools

What's next?

Next Tue, Nov the 7th - two appointments:

- 9:00am, sala verde, lectio magistralis prof. Seth Lloyd (MIT): *Prospects in Quantum Machine Learning*;
- 4:30pm, room I, lecture on the use of IGtools.

・ロト ・ 同ト ・ ヨト ・ ヨト