

SIMCOGEN: A database for collecting the location of similar regions in complete genomes.

Pierre Vincens^{1,2}, Cécile André³, Anne Badel-Chagnon¹, Frédéric Guyon¹ and Serge Hazout¹

(1) Équipe de Bioinformatique Génomique et Moléculaire, INSERM U436, Université Paris 7 - Denis Diderot, case 7113, 2 place Jussieu, 75251 Paris Cedex 05, France

(2) Département de Biologie (FR36), Ecole Normale Supérieure, 46 rue d'Ulm, 75230 Paris Cedex 05, France

(3) Département de Biomathématiques, CHU Pitié Salpêtrière, 91 boulevard de l'Hôpital, 75634 Paris Cedex 13, France

GOAL

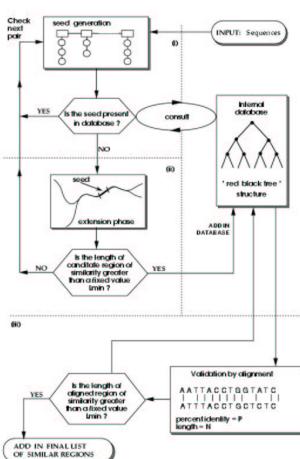
The knowledge of the location of duplicated regions into complete genomes is a major source of information to understand the mechanisms of their evolution [1, 4, 6, 8, 9, 13]. Our primary objective is to create a repository registering all pairs of similar regions found into a lot of complete genomes. Two regions are similar if the identity of their nucleotideic sequences is greater than S_{min} for a length greater than L_{min} .

MATERIAL AND METHODS

DATA OF SIMCOGEN

- Data are registered in a database using a relational model. The main tables include the location of the similar pairs and the informations describing the known regions of analysed chromosomes.
- The search of similar regions has been performed for a lot of chromosomes whose complete sequence is available. Are included prokaryotic (*B. subtilis*, *E. coli*, *M. genitalium*, ...) and eucaryotic organisms (*S. cerevisiae*, ...).
- The nucleotideic sequences and associated informations have been downloaded from:
<ftp://ftp.ncbi.nlm.nih.gov/genbank/genomes>

LOCATING THE SIMILARITIES



- The ASSIRC/DASSIRC [12, 11] approach is used for finding the similarities. The algorithm proceeds in three steps:
- I) Seed generation: A list of seeds is created. A seed is a pair of positions where the sequence has a common motif of fixed size k .
 - II) Seed extension: The bounds of similar regions around the seed are determined. The method computes a 'random walk' [5] on each region from the seed and compares step by step these walks. The pairs of regions are recorded into an internal database by:
 - . Using the Red-Black Tree structure (for rapid access)
 - . Joining adjacent pairs together (to take gaps into account)
 - . Checking for the inclusion of a seed in a previously recorded pair (to reduce the computation).
 - III) Validation of pairs of similar regions: The similarity between two regions is checked using an alignment algorithm [7, 10]

CONCLUSIONS

SIMCOMGEN will make available the location of pairs of similar regions found in complete genomes. Recent works [2] have proved the interest of this repository to study the structure of subtelomeric regions of chromosomes of *Saccharomyces cerevisiae*. However, it is required to complete the development of specific strategies as Mosaic approach [3] to analyse these data.

AVAILABILITY

- ⇒ Interactive: <http://www.biologie.ens.fr/simcogen>
⇒ Flat data: <ftp://ftp.biologie.ens.fr/pub/simcogen>

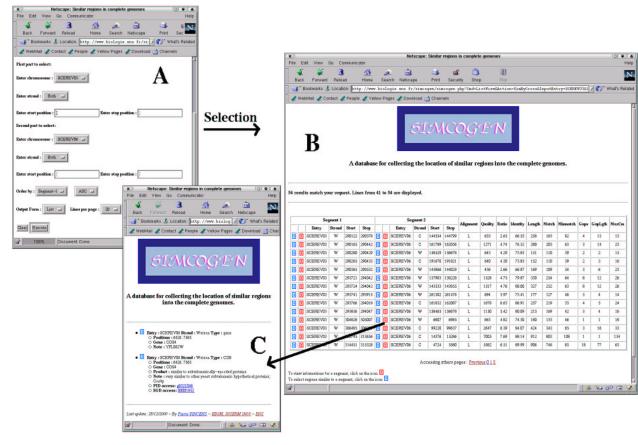
REFERENCES

- [1] G. ACHAZ, E. COISSAC, A. VIARI, and P. NETTER, (2000). Analysis of intrachromosomal duplications in yeast *Saccharomyces cerevisiae*: a possible model for their origin. *Mol Biol Evol*, 17, pp. 1268-75.
- [2] C. ANDRÉ. *Etude des duplications et analyse de l'organisation subtélomérique dans les chromosomes de la levure *Saccharomyces cerevisiae**, PhD thesis, Université Paris VII - Denis Diderot, (2000).
- [3] C. ANDRÉ, P. VINCENS, J. BOISVIEUX, and S. HAZOUT, (2001). Mosaic: Segmenting multi-aligned dna sequences. *Bioinformatics*, 17, pp. 196-197.
- [4] E. COISSAC, E. MAHLER, and P. NETTER, (1997). A comparative study of duplications in bacteria and eukaryotes: the importance of telomeres. *Mol Biol Evol*, 14, pp. 1062-1074.
- [5] M. GATES, (1985). Simpler dna sequence representations. *Nature*, 316, pp. 219-219.
- [6] B. LOFTUS, U. KIM, V. SNEDDON, F. KALUSH, R. BRANDON, J. FUHRMANN, T. MASON, M. CROSBY, M. BARNSTEAD, L. CRONIN, A. D. MAYS, Y. CAO, R. XU, H. KANG, S. MITCHELL, E. EICHLER, P. HARRIS, J. VENTER, and M. ADAMS, (1999). Genome duplications and other features in 12 mb of dna sequence from human chromosome 16p and 16q. *Genomics*, 60, pp. 295-308.
- [7] S. NEEDLEMAN and C. WUNSCH, (1970). A general method applicable to the search for similarities in the amino acid sequence of two proteins. *Journal of Molecular Biology*, 48, pp. 43-453.
- [8] C. SEOGHE and K. WOLFE, (1998). Extent of genomic rearrangement after genome duplication in yeast. *Proceedings of the National Academy of Sciences of the United States of America*, 95, pp. 4447-4452.
- [9] ———, (1999). Updated map of duplicated regions in the yeast genome. *Gene*, 238, pp. 253-61.
- [10] T. SMITH and M. WATERMAN, (1981). Comparison of biosequences. *Advance in Applied Mathematics*, 2, pp. 482-489.
- [11] P. VINCENS, A. BADEL-CHAGNON, F. GUYON, C. ANDRÉ, and S. HAZOUT, Stratégie de recherche des similitudes intrachromosomiques, in JOBIN 2000, G. Caraux, O. Gascuel, and M. F. Sagot, eds., AGRO Montpellier - LIRM, Montpellier, (2000), pp. 383-389.
- [12] P. VINCENS, L. BUFFET, C. ANDRÉ, J. CHEVROLAT, J. BOISVIEUX, and S. HAZOUT, (1998). A strategy for finding regions of similarity in complete genome sequences. *Bioinformatics*, 14, pp. 715-725.
- [13] K. WOLFE and D. SHIELDS, (1997). Molecular evidence for an ancient duplication of the entire yeast genome. *Nature*, 387, pp. 708-713.

RESULTS

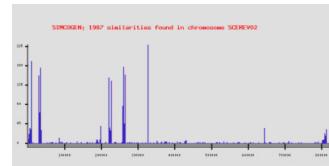
AN EXAMPLE OF INTERACTIVE SESSION

SIMCOGEN can be consulted using a web interface allowing to select and display the pairs of similar regions.

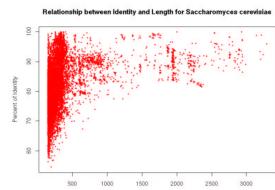


- A: Selection panel to choose the chromosome and the range of positions to consider. Many options allow the user to format the result.
B: Result panel to view the pairs corresponding to the previous selection. For each segment of pairs, the option gives complementary informations and the option searches all similar pairs including this segment.
C: Information panel to display data associated with a segment.

GRAPHIC TOOLS



Distribution of pairs of similar regions along the chromosome II of *S. cerevisiae*



Relationship between the identity and the length of similar pairs found in the genome of *S. cerevisiae*

HETEROGENEITY OF THE DISTRIBUTION OF SIMILAR PAIRS IN A LOT OF STUDIED GENOMES

- The number of intrachromosomal duplications is variable:
 - . high in *Neisseria meningitidis*, *Deinococcus radiodurans chr 1*, *Ureaplasma urealyticum*
 - . low in *Pyrococcus abyssi*, *Pyrococcus horikoshii*, *Treponema pallidum*
- Very large duplications are observed for *Neisseria meningitidis* (3171 bps), *Saccharomyces cerevisiae* (chr 2) (2864 bps), *Treponema pallidum* (1845 bps), *Escherichia coli* (1314 bps).
- *Rickettsia prowazekii* is atypical: the pair of highest quality score has only a similarity of 64.2%.

ID	AC	Name	L	N	lms	pms	fms
BSUBT101	A1009126	<i>Bacillus subtilis</i>	4214814	176	3186	99.43	358
DRADUR01	AE000513	<i>Deinococcus radiodurans chr 1</i>	2648638	1273	2859	99.27	695
DRADUR02	AE001825	<i>Deinococcus radiodurans chr 2</i>	412348	42	600	99.27	695
ECOLIK01	U00096	<i>Escherichia coli</i>	4639221	637	3406	98.76	1314
MYPGEN01	L43967	<i>Mycoplasma genitalium</i>	580074	111	1244	91.68	117
NMENSB01	AE020098	<i>Neisseria meningitidis</i>	2272351	8971	3245	99.48	3171
PABYSS01	A1068636	<i>Pyrococcus abyssi</i>	1765118	51	1382	91.97	873
PYHOR01	Pyro-h	<i>Pyrococcus horikoshii</i>	1738505	75	1592	94.26	623
RPROWA01	AJ235269	<i>Rickettsia prowazekii</i>	1111523	244	328	64.20	26
SCEREV02	NC_001134	<i>Saccharomyces cerevisiae chr 2</i>	813140	1987	2927	99.93	2864
SCEREV04	NC_001136	<i>Saccharomyces cerevisiae chr 4</i>	1531929	205	2756	91.68	305
SCEREV12	NC_001144	<i>Saccharomyces cerevisiae chr 12</i>	1078172	156	2252	88.43	239
SCEREV16	NC_001148	<i>Saccharomyces cerevisiae chr 16</i>	948061	76	1365	98.68	270
TMARIT01	AE000512	<i>Thermotoga maritima</i>	1860725	132	1917	96.92	921
TPALLI01	AE000520	<i>Treponema pallidum</i>	1138011	53	1856	100.00	1845
UUREAL01	AF222894	<i>Ureaplasma urealyticum</i>	751719	1166	764	99.74	759
XFASTI01	AE003849	<i>Xylella fastidiosa</i>	2679306	148	2920	95.64	191

Similarities found for a lot of genomes: ID: SIMCOGEN identifier of the chromosome, AC: Accession number, Name : Current name of the organism, L: Length in nucleotide of the chromosome, N: Number of pairs of intrachromosomal regions of similarity found into the chromosome, lms: Length of aligned sequences for the pair of highest quality score, pms: Percent of identity of aligned sequences for the same pair, fms: Length of the larger identical segment between two similar sequences.