

# Evolutionary Distance in the Yeast Genome

Somayyeh Azizi, Saeed Kaboli, and Atsushi Yagi

**Abstract**—Whole genome duplication (WGD) increased the number of yeast *Saccharomyces cerevisiae* chromosomes from 8 to 16. In spite of retention the number of chromosomes in the genome of this organism after WGD to date, chromosomal rearrangement events have caused an evolutionary distance between current genome and its ancestor. Studies under evolutionary-based approaches on eukaryotic genomes have shown that the rearrangement distance is an approximable problem. In the case of *S. cerevisiae*, we describe that rearrangement distance is accessible by using dedoubled adjacency graph drawn for 55 large paired chromosomal regions originated from WGD. Then, we provide a program extracted from a C program database to draw a dedoubled genome adjacency graph for *S. cerevisiae*. From a bioinformatical perspective, using the duplicated blocks of current genome in *S. cerevisiae*, we infer that genomic organization of eukaryotes has the potential to provide valuable detailed information about their ancestry genome.

**Keywords**—Whole-genome duplication, Evolution, Double-cut-and-join operation, Yeast.

## I. INTRODUCTION

GENOME duplication is a fundamental process in the evolution of species. Within the yeast *Saccharomyces cerevisiae* (*S. cerevisiae*), whole-genome duplication (WGD), is thought to have occurred approximately 100 Ma, in which two identical copies of each chromosome with identical gene content and gene order created [1-5]. Since the WGD to date, organization of the genome of *S. cerevisiae* has been modifying under large chromosomal rearrangements which is critical for the evolution [6]. In other words, during the course of evolution, the genes in the genome can be shuffled around by genome rearrangements that move genes within a chromosome or among chromosomes. For multi-chromosomal genomes, the most common operations of rearrangement are intra-chromosomal translocations, inversions, non-reciprocal telomeric translocations, fusions, fissions [7] and block-interchanges [8] that all were initially introduced by Yancopoulos and colleagues in unique operation to sort two signed genomes which called *double-cut-and-join* (DCJ) [9]. It basically cuts a genome in two places and then joins the resulting four ends in a new way. Recently, Kahn and colleagues constructed an algorithm to compute duplication distance between two signed strings with segmental

duplication. They also extended the computation for the genome with duplication-deletion and duplication-inversion-deletion in the human genome [10]. In [11], genome dedoubling problem was solved with introducing algorithm in the DCJ and reversal models for reconstruction of non-duplicated ancestor of *Drosophila yakuba*. In [12], Chen X and colleagues found how to sort unsigned genome composed of linear and circular chromosome by double-cut-and-join operations. In this paper we calculate the minimum number of rearrangement needed to transform present-day genome of *S. cerevisiae* to its ancestor (DCJ distance) by using theorem 1 mentioned in [11]. To approximate the evolutionary distance between two genomes, we first use updated map of duplicated regions in the yeast genome [2] to draw dedoubled genome adjacency graph. It is well known that the genome of *S. cerevisiae* contains 55 different pairs of syntenic chromosomal regions originated from WGD, consisting of duplicated gene pairs arranged in the same order on two different chromosomes or rarely on a single chromosome [1, 13]. Then, we present a program which is applicable to draw dedoubled genome adjacency graph of *S. cerevisiae* and also any uni- or multi- chromosomal signed genome.

## II. MATERIAL AND METHODS

### A. Whole Duplicated Genome

Each linear chromosome of the genome is composed of duplicated segments and their paralogs which are represented by signed genomic markers. Two copies of a same marker in a genome are called paralogs. Markers are indicated by sign in which the sign shows orientations of markers in chromosomes. Two ends of each linear chromosome indicate telomeres which are represented by the unsigned marker  $\circ$ . If a marker  $x$  is present twice, one of the paralogs is represented by  $\underline{x}$ . A duplicated genome is a genome in which subset of markers are duplicated. For example,  $(\circ \ 2 \ \underline{1} \ -5 \ 6 \ \circ) (\circ \ \underline{2} \ -4 \ \circ) (\circ \ -1 \ \underline{6} \ 3 \ \underline{5} \ \circ)$  is a multi-chromosomal genome composed of four duplicated markers and their paralogs. Whole duplicated genome is a duplicated genome in which all segments are duplicated e.g.,  $(\circ \ 4 \ -5 \ \underline{-5} \ \circ) (\circ \ 2 \ \underline{4-2} \ \circ)$ . Whole duplicated genome in ancestor is a genome such that marker  $x$  and its paralog are in form of pairs  $(x \ \underline{x})$  or  $(\underline{x} \ x)$ . For example,  $(\circ \ 1 \underline{1} \ 2 \ \underline{2} \ 3 \ \underline{3} \underline{4} \ 4 \ 5 \ \underline{5} \ \circ)$ .

### B. Whole Duplicated Genome in the *S. cerevisiae*

A perfectly doubled genome of *S. cerevisiae* included in 16 linear chromosomes rearranged over evolutionary time to construct current genome with cluttered gene order. Present-day rearranged genome of *S. cerevisiae* carrying 55 duplicated pair regions which are surrounded by single-copy genes within

Somayyeh Azizi is now with Department of Applied Physics, Graduate School of Engineering, Osaka University, Osaka, Japan (e-mail: somayyeh\_azizi@ap.eng.osaka-u.ac.jp).

Saeed Kaboli is now with Department of Biotechnology, Graduate School of Engineering, Osaka University, Osaka, Japan. Department of Bioscience, School of Science, Zanzan University, Zanzan, Iran (corresponding author e-mail: kaboli.saeed@bio.eng.osaka-u.ac.jp).

Atsushi Yagi is now with Department of Applied Physics, Graduate School of Engineering, Osaka University, Osaka, Japan (e-mail: yagi@ap.eng.osaka-u.ac.jp).

chromosomes. An overview of location of the 55 duplicated regions on the 16 chromosomes is depicted in Fig. 1.

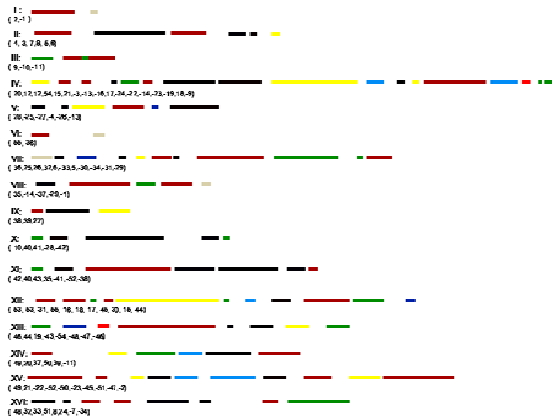


Fig. 1 Locations of the 55 duplicated regions on the 16 chromosomes of the yeast genome

A dedoubled genome of *S. cerevisiae* is a duplicated genome  $G(Sc)$  such that for any duplicated marker  $x$  in  $G(Sc)$ , either  $(x \underline{x})$ , or  $(\underline{x} x)$  is an adjacency of  $G(Sc)$ .

### C. Rearrangement by DCJ Operation

In this paper rearrangement is presented as DCJ operation that acts on two adjacencies  $ab$  and  $cd$  of a genome, transforms them either to the two adjacencies  $ac$  and  $bd$  or two adjacencies  $ad$  and  $bc$ . For example, the following DCJ cuts adjacencies  $(\underline{2} \ 6)$  and  $(-5 \ 2)$  to produce adjacencies  $(\underline{2} \ 2)$  and  $(-5 \ 6)$ . Here, breakpoints are indicated by symbol  $\blacktriangle$ , and new adjacencies are presented by dots.  $(\circ \ \underline{2} \blacktriangle \ 6 \ -3 \ 7 \ \circ)(\circ \ 4 \ -5 \blacktriangle \ 2 \ -3 \ 5 \ \circ) \rightarrow (\circ \ 4 \ -5 \ . \ 6 \ -3 \ 7 \ \underline{2} \ . \ 2 \ -3 \ 5 \ \circ)$

## III. RESULTS

We first study the *S. cerevisiae* genome dedoubling problem under the DCJ model. Then, we present a program which is applicable to draw dedoubled adjacency graph of genomes as well as the genome of *S. cerevisiae*.

### A. Dedoubling of *S. cerevisiae* Genome by DCJ Operation

In order to express the DCJ dedoubled distance,  $d_{dcj}(Sc)$ , between two whole duplicated genomes, we use a graph called dedoubled adjacency graph of *S. cerevisiae*.

### B. Dedoubled Adjacency Graph in *S. cerevisiae*

The dedoubled adjacency graph of *S. cerevisiae* denoted by  $A(Sc)$ , is the graph consists of two edges for each marker of  $G(Sc)$ , one edge between the vertices  $(x \ u)$  and  $(v \ \underline{x})$ , another between the vertices  $(y \ x)$  and  $(\underline{x} \ z)$ . For example, marker -13 into the  $G(Sc)$  contains two edges between the vertices  $(-13 \ -16)$ ,  $(-26 \ -13)$  and  $(3 \ -13)$ ,  $(-13 \ \circ)$  respectively (see Fig. 2).

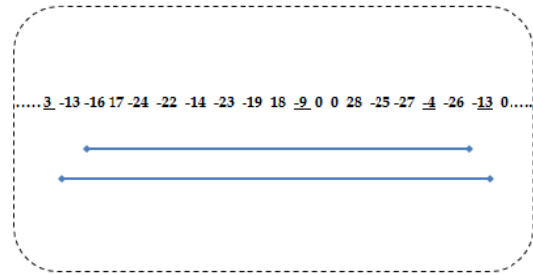


Fig. 2 Vertices and edges for marker -13 in  $A(Sc)$

Vertices are the adjacencies of  $G(Sc)$ . Note that all vertices in  $A(Sc)$  have degree one or two. Thus, the connected components of  $A(Sc)$  are only paths and cycles. These paths and cycles are called elements of  $A(Sc)$ . If an element contains twice the couple  $(x, \underline{x})$  so is called a duplicated element of  $A(Sc)$ , otherwise is called non-duplicated, namely if an element of  $A(Sc)$  contains no couple  $(x, \underline{x})$  twice, then it is non-duplicate elements of this graph. If two edges of each marker belong to two different elements of  $A(Sc)$ , so both elements contain couple  $(x, \underline{x})$  and then these two elements intersect. Otherwise are called independent.

Post-WGD information of *S. cerevisiae* (<http://acer.gen.tcd.ie/~khwolfe/yeast>) helped us to draw dedoubled adjacency graph of the present-day genome of budding yeast with 55 pairs of chromosomal markers as depicted in Fig. 3.

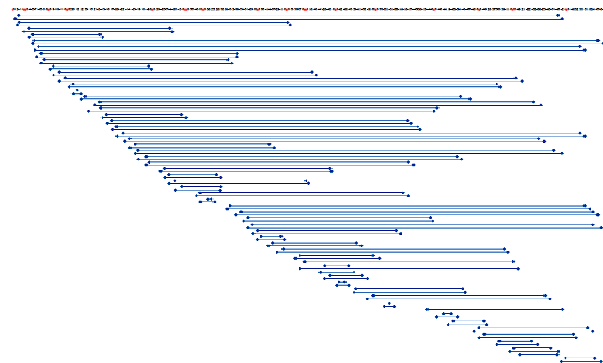


Fig. 3 The dedoubled adjacency graph of *Saccharomyces cerevisiae* genome

Based on the drawn graph, 16 elements were resulted in which 13 elements are paths while remaining 3 elements are cycles. All constructed cycles, i.e.  $C_1; \{(12, \underline{12})\}$ ,  $C_2; \{(53, \underline{53})\}$  and  $C_3; \{(14 \ -37) \ (37 \ 50) \ (-50 \ -23)\} = (-37, 37) \ (50, \underline{50})\}$  are non-duplicated cycles. None of these non-duplicated cycles have common couple  $(x, \underline{x})$  together. Thus, they are non-duplicated pairwise independent cycles.

According to theorem 1 of [11], for *S. cerevisiae* $_{dcj}(Sc)$  would be equal with  $n - C_i$ . In this equation, let  $n$  be the number of couple of paralogous marker in  $G(Sc)$  and let  $C_i$  be the maximum size of subset of non-duplicated pairwise independent cycles in  $A(Sc)$ . Thus, according to Fig. 3

Fig. 4 Program for drawing dedoubled adjacency graph of *Saccharomyces cerevisiae* genome

#### IV. CONCLUSION

From a bioinformatical perspective, using the duplicated blocks in the current genome of *S. cerevisiae*, we have inferred that how to generalize the dedoubled adjacency graph for whole genome of this organism which is valuable for deep understanding of rearrangement events since WGD to date. In this work, we have shown that the number of non-duplicated pairwise independent cycle of current genome negatively effects on its distance with ancestor genome.

- [1] Wolfe KH, Shields DC, "Molecular evidence for an ancient duplication of the entire yeast genome", *Nature* 387, 708±713, 1997.
- [2] Seoighe C, Wolfe KH, "Updated map of duplicated regions in the yeast genome", *Gene* 238, 253–261, 1999.

- [3] Kellis M, Birren BW, Lander ES, "Proof and evolutionary analysis of ancient genome duplication in the yeast *Saccharomyces cerevisiae*", *Nature* 428: 617-624, 2004.
- [4] Vitkup D, Kharchenko P, Wagner A, "Influence of metabolic network structure and function on enzyme evolution", *Genome Biol* 7: R39, 2006.
- [5] Conant GC, Wolfe KH, "Turning a hobby into a job: how duplicated genes find new functions", *Nat Rev Genet* 9: 938-950, 2008.
- [6] Gordon JL, Byrne KP, Wolfe KH, "Additions, losses, and rearrangements on the evolutionary route from a reconstructed ancestor to the modern *Saccharomyces cerevisiae* genome", *Plos Genet* 5(5) e1000485, 2009.
- [7] Gordon JL, Byrne KP, Wolfe KH, "Mechanisms of chromosome number evolution in yeast", *Plos Genet* 7(7) e1002190, 2011.
- [8] Christie, D.A., "Sorting by block-interchanges", *Inf. Process. Lett.* 60, 165-169, 1996.
- [9] Yacopoulos S, Attie O, Friedberg R, "Efficient sorting of genomic permutations by translocation, inversion and block interchange", *Bioinformatics* 21(16): 3340-3346, 2005.
- [10] Kahn CL, Mozes S, Raphael B J, "Efficient algorithms for analyzing segmental duplications with deletions and inversion in genomes", *Algorithms for Molecular Biology* 5: 11, 2010.
- [11] Thomas A, Varré JS, Ouangraoua A, "Genome dedoubling by DCJ and reversal", *BMC bioinformatics* 12(Suppl 9):S20, 2011.
- [12] Chen X, Sun R, Yu J, "Approximating the double-cut-and-join distance between unsigned genomes", *BMC Bioinformatics* 12 (Suppl 9): S17, 2011.
- [13] Mewes HW, Albermann K, BaÉhr M, Frishman D, Gleissner A et al., "Overview of the yeast genome", *Nature* 387, Suppl., 7-65, 1997.

**Somayyeh Azizi** was born in Zanjan, Iran. She received the B. D. in the Department of Computer Engineering from Zanjan Azad University in 2010. She is currently in the five years doctor's course in the Department of Applied Physics, Osaka University. She is interested in Bioinformatics.

**Saeed Kaboli** was born in Zanjan, Iran on August 23, 1978. He received the M.S. degree in the Department of Medical Mycology and Parasitology from Mazandaran University of Medical Sciences, Sari, Iran in 2008. He is currently in the doctor's course in the Department of Biotechnology, Graduate School of Engineering, Osaka University, Osaka, Japan. He is interested in Fungal Genetics and Evolutionary Genomics of Industrially and Medically Important Fungi.

**Atsushi Yagi** was born in Niigata, Japan in 1951. He received the Ph.D. degree in Division of Natural Science from Osaka University in 1982. He is currently a full Professor of Department of Applied Physics, Graduate School of Engineering, Osaka University, Osaka, Japan. His research interests include Evolution Equations, Self-Organization Model and its application in Physics, Chemistry and Biology.