

Complete nucleotide sequence of *Saccharomyces cerevisiae* chromosome X

F. Galibert^{1,2}, D. Alexandraki³, A. Baur⁴, E. Boles⁴, N. Chalwatzis⁴, J.-C. Chuat¹, F. Coster⁵, C. Cziepluch⁶, M. De Haan⁷, H. Domdey⁸, P. Durand⁹, K. D. Entian¹⁰, M. Gatius¹, A. Goffeau⁵, L. A. Grivell⁷, A. Hennemann¹⁰, C. J. Herbert¹¹, K. Heumann¹², F. Hilger⁹, C. P. Hollenberg¹³, M.-E. Huang¹, C. Jacq¹⁴, J.-C. Jauniaux⁶, C. Katsoulou³, L. Kirchrath¹³, K. Kleine¹², E. Kordes⁶, P. Kötter¹⁰, S. Liebl¹², E. J. Louis¹⁵, V. Manus¹, H. W. Mewes¹², T. Miosga⁴, B. Obermaier¹⁶, J. Perea¹⁴, T. Pohl¹⁷, D. Portetelle⁹, A. Pujol⁸, B. Purnelle⁹, M. Ramezani Rad¹³, S. W. Rasmussen¹⁸, M. Rose¹⁰, R. Rossau¹⁹, I. Schaaff-Gerstenschläger⁴, P. H. M. Smits⁷, T. Scarcez¹⁹, N. Soriano¹, D. Tovan¹⁴, M. Tzermia⁷, A. Van Broekhoven¹⁹, M. Vandenbol⁹, H. Wedler²⁰, D. Von Wettstein¹⁸, R. Wambutt²⁰, M. Zagulski^{11,21}, A. Zöllner¹² and L. Karpfinger-Hartl¹²

¹UPR 41 CNRS Recombinations Génétiques, Faculté de Médecine, 2 avenue de Professeur Léon Bernard, F-35043 Rennes Cedex, France, ²Foundation for Research and Technology Hellas, Institute of Molecular Biology and Biotechnology, PO Box 1527, Heraklion, GR-71110 Crete, Greece, ³Institut für Mikrobiologie und Genetik, Technische Hochschule Darmstadt, Schriesheimstrasse 10, D-64287 Darmstadt, Germany, ⁴Unité de Biochimie Physiologique, Université Catholique de Louvain, Place Croix du Sud 2, Bâtiment 20, B-1348 Louvain-La-Neuve, Belgium, ⁵Tumorsvirologie Abteilung 0610 and Virologie Appliquée à l'Oncologie Unité INSERM U375, Deutsches Krebsforschungszentrum, D-69120 Heidelberg, Germany, ⁶University of Amsterdam, Section for Molecular Biology, Kruislaan 318, NL-1098 SM Amsterdam, The Netherlands, ⁷Genzentrum, Institut für Biochemie, Würmstrasse 221, D-81373 München, Germany, ⁸Unité de Microbiologie, Faculté des Sciences Agronomiques de Gembloux, avenue Maréchal Juin 5, B-5000 Gembloux, Belgium, ⁹Institut für Mikrobiologie, J.W. Goethe-Universität Frankfurt, Marie-Curie-Strasse 9, Geb. N250, D-60439 Frankfurt/Main, Germany, ¹⁰UPR 2420 CNRS Centre de Génétique Moléculaire, Bâtiment 26, Avenue de la Terrasse, F-91198 Gif-sur-Yvette cedex, France, ¹¹MIPS am Max-Planck-Institut für Biochemie, D-82152 Martinsried bei München, Germany, ¹²Institut für Mikrobiologie der Heinrich-Heine-Universität Düsseldorf, Geb. 26.12, Universitätsstrasse 1, D-40225 Düsseldorf, Germany, ¹³URA 1302 CNRS Génétique Moléculaire, Ecole Normale Supérieure, 46 rue d'Ulm, F-75230 Paris Cedex 05, France, ¹⁴Yeast Genetics, Institute of Molecular Medicine, John Radcliffe Hospital, Headington, Oxford OX3 9DU, UK, ¹⁵GAZC GmbH, Gesellschaft für Analytische Technik und Consulting, Fritz-Arnold-Strasse 23, D-78467 Konstanz, Germany, ¹⁶Carlsberg Laboratory, Department of Physiology, Gamle Carlsberg vej 10, Valby, DK-2300 Copenhagen, Denmark, ¹⁷Imogenetics, Industriepark Zwijpsande 7, Box 4, B-9052 Ghent, Belgium and ¹⁸ADON GmbH, Gesellschaft für molekularbiologische Technologie mbH, Glienicke Weg 185, D-12489 Berlin, Germany

¹⁹Present address: MediGene GmbH, Lochthamer Strasse 11, D-82152 Martinsried bei München, Germany

²⁰Present address: Institute of Biochemistry and Biophysics, 5a Pawinskiego St., 02-106 Warsaw, Poland

²¹Corresponding author

The complete nucleotide sequence of *Saccharomyces cerevisiae* chromosome X (745 442 bp) reveals a total of 379 open reading frames (ORFs), the coding region covering ~75% of the entire sequence. One hundred and eighteen ORFs (31%) correspond to genes previously identified in *S. cerevisiae*. All other ORFs represent novel putative yeast genes, whose function will have to be determined experimentally. However, 57 of the latter subset (another 15% of the total) encode proteins that show significant analogy to proteins of known function from yeast or other organisms. The remaining ORFs, exhibiting no significant similarity to any known sequence, amount to 54% of the total. General features of chromosome X are also reported, with emphasis on the nucleotide frequency distribution in the environment of the ATG and stop codons, the possible coding capacity of at least some of the small ORFs (<100 codons) and the significance of 46 non-canonical or unpaired nucleotides in the stems of some of the 24 tRNA genes recognized on this chromosome.

Keywords: chromosome X/gene duplication/open reading frame/*Saccharomyces cerevisiae*/tRNA

Introduction

The traditional methods of genetic analysis involve tracing modified phenotypes back to genotypic alterations. The limit of this approach is an imperceptible modification of the phenotype. The international yeast genome systematic sequencing programme launched in 1989 by the European Communities, aiming at establishing the complete genetic information of bakers' yeast, *Saccharomyces cerevisiae*, has demonstrated the limitations of classical genetics. The pilot sequencing of chromosome III (Oliver *et al.*, 1992) has demonstrated that disruption of a large number of the newly revealed open reading frames (ORFs) does not result in any phenotypic alteration. Subsequent systematic sequencing of seven more chromosomes (Barrell *et al.*, 1994; Dietrich *et al.*, 1994; Dajon *et al.*, 1994; Feldmann *et al.*, 1994; Johnston *et al.*, 1994; Bussey *et al.*, 1995; Murakami *et al.*, 1995) has confirmed that a large proportion of the novel genes cannot be assigned any known function, while on the other hand a large number of proteins unrelated to database entries are being discovered. Last but not least, it stems from numerous cytological studies of chromosome behaviour during the vegetative and meiotic cell cycle that a chromosome is more than its mere genetic content. By making available the complete

Table I. Estimated overall accuracy of chromosome X sequence

	Total bp verified	Number of modified nt ^a			Error rate (%)
		M	G	T	
Overlap between regions	46 455	11	13	24	0.52
Resequenced regions ^b	~50 000	10	7	17	0.34

^aM, mismatch; G, gap; T, total mismatches plus gaps.

^bOccasional overlaps between verification clone sequences were excluded from the calculations.

DNA sequence of a chromosome, parameters not entirely confined to its role as carrier of genetic information may be exposed for analysis. A survey of a new object is thus provided, even though all the topological implications of the results cannot be fully grasped at the present stage and must await at least the completion of the yeast genome enterprise. This paper describes the DNA sequence of chromosome X.

Results

Assembly of the sequence

The sequence was determined from a set of 26 partially overlapping cosmids selected on the basis of an *EcoRI* map based on a cosmid contig of chromosome X (Huang *et al.*, 1994a). These cosmids were distributed within a consortium of 15 contractors. The telomeres were independently isolated and sequenced. While the left-telomere-containing clone was found to overlap with the left terminal cosmid of the chromosome, this was not so at the other end, where no overlap was detected between the right-most cosmid and a right-telomere-containing clone 9.0 kb in size. The missing portion (a few kb) was PCR-amplified from a yeast S288C genomic DNA template using primers designed from sequences flanking the gap. When all bases had been determined by each contractor and each sequencing strategy had been approved by the DNA coordinator, ensuring that the sequence had been independently determined on each strand with sufficient overlap between all the subclones, the sequences were considered as final and entered into the MIPS data library for assembly. Partial sequences of chromosome X have been published independently by some of the authors of this work (Huang *et al.*, 1994b, 1995; Miesga *et al.*, 1994a,b,c, 1995; Purnelle *et al.*, 1994; Vandenbol *et al.*, 1994, 1995; Rasmussen, 1995; Zagalski *et al.*, 1995).

Verification of the sequence

Quality controls were performed concomitantly with sequence assembly. The aim of the project was to keep the error rate as low as possible, with a target $<10^{-4}$. These procedures were employed to track down errors, including checking sequencing strategy by the coordinator, matching overlapping portions sequenced by independent contractors and finally random resequencing (see Materials and methods for details). The results of the last two procedures are shown in Table I. From these data, the error rate of the yeast chromosome X sequence presented in this paper can be estimated to be 0.4%, a value of the same order as that reported in similar studies.

General organization of chromosome X

Analysis of the entire nucleotide sequence of chromosome X (745 442 bp) confirms the general features of chromosome organization observed in other systematically sequenced yeast chromosomes. The coding region occupies 74.04% of the sequence, 36.59% and 37.45% on the Watson and Crick strand, respectively.

The average base composition is 38.9% G+C. As expected, the coding regions have a higher than average G+C content (40.2%) than the non-coding (35.6%). The distribution of dinucleotide frequencies over the whole chromosome is the same in the coding and the non-coding regions of either strand. The deviations of the frequencies of complementary dinucleotide pairs tend to occur in the same direction. In contrast to what was reported for chromosomes XI and II, the homopurine pairs do not seem to be in excess in the coding region of either strand (Figure 1). Some compositional periodicity has been noted, at least in the case of chromosomes XI and II, with waves of G+C-rich regions correlating with waves of high gene density. By using the same algorithm, a similar G+C pattern emerges with chromosome X, especially in the right-hand part of the chromosome. This pattern correlates rather well with the gene density plot, as illustrated by the two deep depressions around 200 kb and 470 kb in Figure 2.

Telomeres and centromere

The telomere regions of chromosome X are similar to the other sequenced yeast telomeres. Adjacent to the C_{1-3} A repeat at the left telomere are a Y' element (coordinates 61–6931) and the core X element (7305–7767) shared by most if not all yeast telomeres (Louis *et al.*, 1994; Pryde *et al.*, 1995). However, the X–Y' junction does not contain the usual subtelomeric repeats STR-D, STR-C, STR-B and STR-A, but instead has (6998–7224) part of a copy (Louis and Haber, 1991) of the fourth intron of cytochrome *b* encoded by mitochondrial DNA (Delehotde *et al.*, 1989). A copy of *h14* is also found at the left telomere of chromosome IX (Louis and Haber, 1991; Barrell *et al.*, 1994). In fact, the left ends of chromosomes IX and X share a large, nearly identical block of sequence similarity spanning >21 kb. The right telomere of chromosome X is more conventional, with a core X element (744 593–745 052) and the STR-D, STR-C, STR-B and STR-A elements adjacent to the TG_{1-3} repeats (745 357–end). The core X elements of both ends contain the *ARSJ* consensus and the *Abf1p* binding site found in most core Xs. These elements that are shared by most ends may have functional significance. The right telomere region is analogous to several other sequenced telomeres (III right

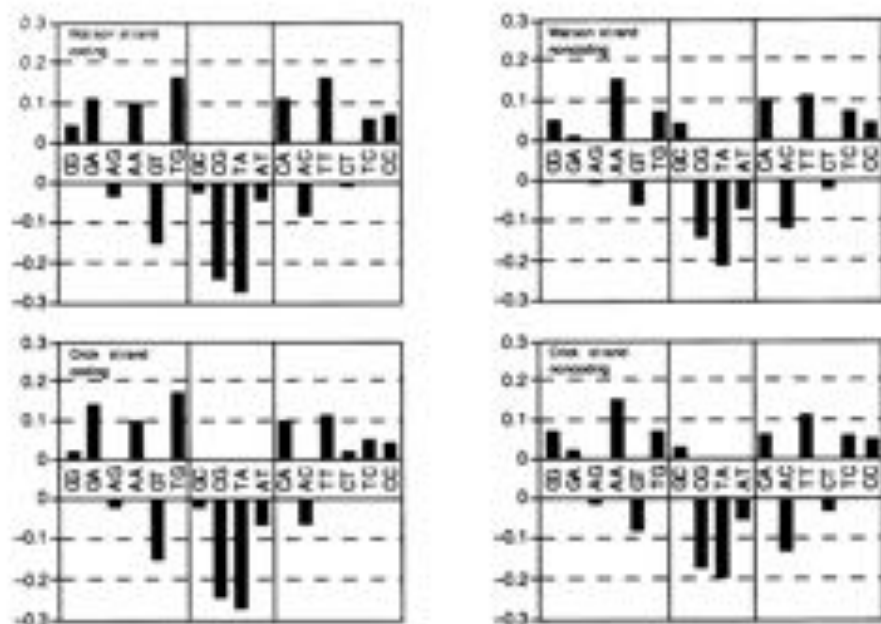


Fig. 1. Distribution of dinucleotide frequencies in the coding and non-coding regions of the two strands of chromosome X. Vertical bars show relative deviations [i.e. (observed-expected)/expected]. Expected frequencies are calculated from mononucleotide frequencies. Complementary pairs are arranged as mirror images. The four self-complementary pairs are placed in the central part.

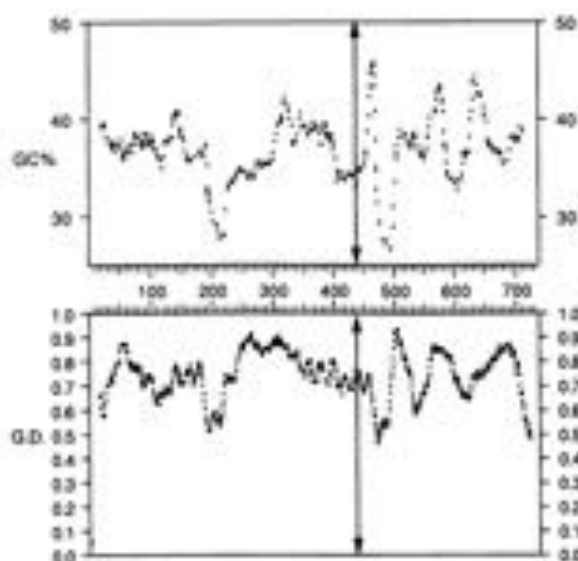


Fig. 2. Compositional variation and gene density distribution along chromosome X. Top: compositional variation calculated as in Dujon *et al.* (1994). Each point represents the average G+C composition calculated from the third base of each codon. Bottom: gene density expressed as the fraction of nucleotides within ORFs in sliding windows of 30 kb. The position of the centromere is indicated by an arrow.

and left, V right and left, VI left, VIII right and left, IX right, XI left) over the last 3–4 kb.

The centromere of chromosome X of strain R95-4A, a derivative of S288C, was isolated by Hieter *et al.* (1985) by selection of yeast DNA fragments capable of suppressing lethality of the *SUP11* gene in high copy number. Comparison of this sequence with that reported in the present

paper shows complete identity and enables location of the chromosome X centromere at positions 435 996–436 112. *CEN10* conforms to the consensus structure established for other centromeres.

ORFs and their predicted protein products

By definition, an ORF is considered from its first in-phase ATG codon. Only those ORFs containing at least 99 contiguous sense codons following an ATG, and not entirely contained within a longer ORF in a different reading frame or on the other DNA strand, have been retained for further analysis. The special case of ORFs shorter than 100 codons is described below. A total of 379 ORFs were recorded in the entire chromosome X using this principle (Table II), leaving aside the retroposons, i.e. a density of one ORF/1967 bp. Twelve of these ORFs are interrupted by introns. Table II includes 39 partially overlapping ORFs. Ten are on the same DNA strand, all others being antiparallel overlaps. Informatic and statistical analysis revealed that ORFs both shorter than 150 codons and with a codon adaptation index (CAI) (Sharp and Li, 1987) <0.11 may correspond to randomly occurring ORFs rather than to real genes (Dujon *et al.*, 1994). If these criteria are applied to the ORFs identified in chromosome X, 23 of the 379 ORFs are questionable genes. Thirteen of these belong to the set of partially overlapping ORFs. However, three genes of known function (*HAP17*, *STE18* and *RPL46*) fall into this category as well, making the border between ORF and gene even more elusive. Taking into account the physical position and ATG environment may help tell which ORFs are genes.

Comparison of the nucleotide sequence and of the predicted protein products with public database entries reveals that 118 ORFs (31%) correspond to genes previously identified in *S.cerevisiae*. All other ORFs represent

Table II. List of ORFs longer than 99 sense codons, known genes and other genetic elements of chromosome X

Nomenclature	Size (aa)	Coordinates	Locus	CAI	Fasta score	Description (nature of element, function or similarity of product)/Comment	
Working Official		1 80				left telomere sequence (complement TG _n -)	
		81 8031				T' element	
X020	YX.122c	3508 869 4130		0.13		probable nucleotide-binding protein, TMM 1+1 (intron from 4582 to 4995)	E
		4998 7224				copy of part of 641 intron from cytochrome <i>b</i> gene (mitochondrial DNA)	
		7305 7307				core X element	
X028	YX.122c	520 8779 9138		0.65	536 (538)	similar to PGL1 protein (PIR: S48516)	B
X0213	YX.122c	1549 11475 16121		0.18	5526 (7778)	similar to carboxypeptidase Y-sorting protein PEP1 (PIR: S23326), TMM 3+1	B
X0218	YX.122c	589 16770 18536		0.25	2459 (3084)	similar to α -glucosidase MAL35 (PIR: S46183), TMM 1+0	B
X0220	YX.122c	130 18240 18892		0.10		hypothetical protein, TMM 2+1	E
X0222	YX.121a	567 19976 21197		0.17	2913 (2955)	similar to lysine transport protein LGT3 (PIR: A51571), TMM 8+1	B
X0226	YX.121a	196 21973 22580		0.18	453 (963)	similar to galactonide O-acetyltransferase (SW: P07864), TMM 1+0	C
X0228	YX.121a	196 23133 23726		0.12		hypothetical protein	F
X0228	YX.121a	581 24344 26086		0.23	2129 (3085)	similar to α -glucosidase (PIR: S45177), TMM 1+0	C
X0231	YX.121a	109 26415 26771		0.16		hypothetical protein, ?	F
X0232	YX.121a	589 26887 28593		0.20	2953 (3023)	probable lysine transport protein RXT6 (PIR: S45176), TMM 11+1	B
X0234	YX.121a	331 32363 33153		0.14		hypothetical protein	F
X0236	YX.121a	799 33853 36249		0.18	1640 (4157)	similar to <i>L</i> -proline OPG (PIR: S45181), TMM 10+1	D
X0238	YX.121a	147 36760 37290		0.10		hypothetical protein, ?	F
X0240	YX.121a	271 36919 37731	CRT1	0.09		CRT1 protein (PIR: S27422)	A
X0242	YX.208a	654 38865 39966	CBP1	0.15		CBP1 protein (PIR: S05826)	A
X0310	YX.208c	329 40197 41183	NLC1	0.14		nuclear NLC1 precursor, mitochondrial (PIR: S05888)	A
X0312	YX.207c	2044 41392 47433		0.14		hypothetical protein, TMM 4+1	E
X0316	YX.206c	758 42962 49933		0.15		hypothetical protein, TMM 1+1	E
X0318	YX.205c	187 50832 51192		0.14		hypothetical protein	F
X0320	YX.206c	645 51216 51630		0.16		hypothetical protein	F
X0322	YX.207a	280 53360 54179	SPPY1	0.14		pre-mRNA splicing factor SPPY1 (PIR: S23553)	A
X0323	YX.202c	115 53945 54289		0.12		hypothetical protein, TMM 1+1	E
X0325	YX.201a	599 54378 56174		0.15		hypothetical protein	F
X0327	YX.208c	789 56446 58032		0.22	2130 (3762)	similar to mitochondrial acetate hydratase (CB: U37788)	C
X0330		58999 59171				RNA ^{7m}	
X0332		59471 59782				I element	
X0334	YX.199c	108 59837 60180		0.09		hypothetical protein, ?	F
X0336	YX.199c	881 60842 61884		0.18	2799 (4118)	similar to YCR037c (PIR: S46655), TMM 13+1	C
X0340	YX.197a	1254 63865 67564		0.14	535 (617)	probable ubiquitin-carboxyl terminal hydrolase (SW: P29125)	D
X0343	YX.199c	590 67831 68780		0.13	924 (1753)	similar to steroid isomerase S4384 (PIR: S46636), TMM 5+0	C
X0345	YX.195c	233 69242 69980		0.11		hypothetical protein, TMM 2+0	F
X0347	YX.194a	513 69336 70874	CDC9	0.13		cell division control protein CDC9 (PIR: S46640)	A
X0349	YX.195a	402 71364 72569		0.10	447 (2131)	similar to SLY41 protein (PIR: S46641), TMM 6+1	D
X0351	YX.192c	234 72711 73412		0.16		hypothetical protein, TMM 2+0	E
X0353	YX.191a	136 73785 74696	CRV2	0.39		ribosomal protein S146B (intron from 73795 to 74202) (PIR: S46642)	A
X0355	YX.190c	130 74911 75308	AP32A	0.30		ribosomal protein S15d (PIR: A23802)	A
X0360	YX.189a	51 75931 76469	RP24b	0.92		ribosomal protein L39c (intron from 75937 to 76172) (EMBL: X01963)	B
X0403	YX.188c	102 76265 76598		0.15		hypothetical protein	F
X0408	YX.187c	819 76804 79260	SWK1	0.13		protein kinase SWK1 (PIR: S46900), TMM 1+0	A
X0409	YX.186a	586 80152 81909		0.16	1039 (3006)	similar to TTP1 protein (PIR: S45576), TMM 2+0	C
X0413	YX.185c	293 82085 82573		0.11		hypothetical protein	F
X0420	YX.184a	123 83445 83813		0.08		hypothetical protein, ?	F
X0423	YX.183a	422 84085 85130		0.18		hypothetical protein, TMM 1+0	E
X0430	YX.182c	103 85435 85749		0.08		hypothetical protein, TMM 1+0, ?	E
X0433	YX.181a	611 85627 87489		0.11	443 (2905)	hypothetical protein, similar to J1575, TMM 1+1	F
X0486	YX.180c	525 87583 88557	ATP12	0.12		ATP12 protein precursor (PIR: A39736)	A
X0488	YX.179a	109 88784 89110		0.15		hypothetical protein	F
X0490	YX.178c	196 89282 89849		0.17		hypothetical protein, TMM 1+0	E
X0493	YX.177a	164 90782 91651		0.68	823 (827)	ribosomal protein L17c (intron from 91093 to 91405) (PIR: S38612)	B
X0495	YX.176c	825 92052 94526	SW1	0.15		transcription factor SW1 (PIR: S26786)	A
X0502	YX.175a	170 94045 94554		0.12		hypothetical protein, TMM 3+0	E
X0504	YX.174a	276 95088 95915	KRE9	0.16		secretory pathway protein KRE9 precursor (PIR: S23891), TMM 1+0	A
X0506	YX.173c	122 96180 96325	RAJ1	0.14		replication factor A chain 3 (PIR: C37281)	A
X0510	YX.172a	411 97729 99456	CPH1			Gly-X carboxypeptidase precursor (PIR: S14605)	A
X0512	YX.171c	396 99699 100886		0.22	478 (1923)	hypothetical protein, similar to YBR162C (PIR: S48655), TMM 2+0	D
X0514	YX.170c	183 101045 101493		0.13		hypothetical protein, TMM 2+0	E
X0517	YX.169a	122 102890 103455		0.15		hypothetical protein, TMM 2+0	E
X0520	YX.168c	713 103223 104419		0.14	258 (1093)	similar to ribonuclease ALL1 (see Intron model) (PIR: A44284)	D
X0523	YX.167a	282 105005 106080	PPP1			ferrous pyrophosphate synthase (SW: A34441), TMM 1+1	A
X0526	YX.166a	94 106425 106706		0.23	QC85	ubiquinol-cytochrome <i>c</i> reductase subunit VII (PIR: S48136)	A
X0531	YX.165c	825 106888 109452	HAL3	0.13		HAL3 protein (PIR: S48240)	A
X0541	YX.164c	387 109960 111150	RAJ1	0.14		protein kinase, cAMP-dependent, catalytic chain 1 (PIR: A27070)	A
X0544	YX.163c	555 111662 113526		0.08		hypothetical protein, TMM 11+1	E
X0549	YX.162c	482 114177 115622		0.14		hypothetical protein	F

Table II. Continued

Nomenclature	Size (aa)	Coordinates	Locus	CAF Fasta score	Description (nature of element, function or similarity of product)/Comment	
Working Official						
Y050		115032-116003			rDNA ²⁰	
Y052	YR104a	180-117238-117777		0.09	hypothetical protein, TMM 0+1	E
Y053	YR105a	180-118260-118819		0.15 326 (1751)	similar to PIR1 protein (ye X1) (PIR: S3560)	C
Y058	YR109a	319-120443-121372		0.47 377 (1162)	similar to PIR2 protein (ye X1) (PIR: S3561)	C
Y061	YR109c	227-121968-122644		0.39 321 (976)	similar to PIR2 protein (ye X1) (PIR: S3561)	C
Y065	YR107c	830-121035-126024	Y065	0.13	LacIa arrest protein YAR1 (SW: S13341)	A
Y070	YR109c	687-126599-126649		0.13	hypothetical protein, TMM 0+1	E
Y073	YR105c	452-128985-130340	FBP20	0.14	fructose-2,6-bisphosphate 2-phosphatase (PIR: A42569)	A
Y080	YR104c	944-130801-130617	VPS37	0.15	vacuolar protein-sorting protein VPS35 (PIR: S70295)	A
Y0810	YR105c	555-134832-135046	Y0810	0.14	myo-inositol 1-phosphate synthase (PIR: A39902), TMM 2+1	A
Y0828	YR102c	119-135871-136227		0.07	hypothetical protein, TMM 0+0, 7	E
Y0830	YR101c	133-136872-136470		0.36	hypothetical protein, TMM 2+0	E
Y0832	YR100a	309-136820-137119		0.09	hypothetical protein, TMM 0+0, 7	E
Y0834	YR109a	663-137076-139064		0.34 298 (1276)	hypothetical protein, similar to YDR032.06c (GB: S51858), TMM 1+0	E
Y0835		139408-139647	SNR190		SnR 190 small nuclear RNA	
Y0836		139263-140360	SNR129		SnR 129 small nuclear RNA	
Y0837	YR108a	233-140134-140302		0.20	hypothetical protein	F
Y0839	YR107c	382-141119-142264		0.15	hypothetical protein	F
Y0842	YR106a	469-142989-144395		0.11	hypothetical protein, TMM 1+0	E
Y0844	YR109a	294-144897-145736		0.22	hypothetical protein	F
Y0846	YR104a	304-146936-148367		0.07	hypothetical protein, 7	F
Y0848	YR107a	134-148798-147271	MM17	0.14	mitochondrial inner membrane protein MM17 (PIR: S46257), TMM 1+1	A
Y0850	YR102c	130-147519-147908		0.06	hypothetical protein, TMM 1+1, 7	E
Y0852	YR101c	807-143863-150807	Y0852	0.12	protein kinase YAK1 (PIR: A32582), TMM 1+0	A
Y0854	YR100a	224-150658-151720	RPB4	0.14	DNA-directed RNA polymerase II chain RPB4 (PIR: A32646)	A
Y0857	YR109c	428-151433-152096	Y18Y	0.14	YLR1 protein (PIR: S26426), TMM 1+0	A
Y0860	YR106c	395-151204-151388	Y072	0.79	translation initiation factor eIF-4A (GB: X12814)	A
Y0863	YR107c	380-154663-155820		0.14 445 (1978)	hypothetical protein, similar to YKR059a (PIR: S38154)	D
Y0864	YR109c	87-156283-156970		0.06	ribosomal protein S21a (accession from U6487 to U6494)	B
Y0866	YR105a	105-157374-157686		0.14	hypothetical protein	F
Y0871	YR104a	409-157885-159111		0.10 1298 (2132)	hypothetical protein, similar to YKR075c (PIR: S38127), TMM 4+1	F
Y0873	YR107a	314-160316-161237	MRS1	0.08	splicing protein MRS1, mitochondrial (PIR: S01267)	A
Y0878	YR102c	780-161811-162889		0.12	hypothetical protein, TMM 1+1	E
Y0882	YR101c	336-163978-165045		0.12	hypothetical protein	F
Y0886	YR109c	224-165423-172664	URA2	0.29	pyrimidine synthesis protein URA2 (PIR: S07367), TMM 1+1	A
Y0889	YR106a	105-171926-172829		0.06	hypothetical protein, access from U72802 to U72805, 7	F
Y0893	YR102a	1235-173299-173963	Y0893	0.14	potassium transport protein, high-affinity (PIR: S05609), TMM 8+1	A
Y0899	YR107c	668-177797-178689	PR12	0.14	polymyxin B resistance protein kinase (PIR: A32794)	A
Y0902	YR107c	640-181999-183934	SPY30	0.12	regulatory protein SPY30 (PIR: S47865)	A
Y0906	YR102a	307-184399-185339		0.12 309 (1519)	hypothetical protein, similar to L08383 (GB: U09102)	F
Y0910	YR102c	363-185229-186377		0.14	hypothetical protein	F
Y0914	YR104c	172-186828-187343		0.16	hypothetical protein	F
Y0918	YR102c	478-187796-189139		0.15	hypothetical protein	F
Y0923	YR102a	178-188415-189939		0.21	hypothetical protein	F
Y0924	YR101c	238-188876-190799	RPY2	0.30	ribulose-5-phosphate 3-epimerase (GB: 63771)	A
Y0934	YR102a	297-190721-191043		0.14	hypothetical protein, TMM 1+1	E
Y0938	YR109c	107-191274-191544		0.13	hypothetical protein, TMM 1+0	E
Y0942	YR109a	219-191336-191994		0.09	hypothetical protein, TMM 1+1	E
Y0944	YR107a	311-192230-193162		0.19	hypothetical protein, TMM 2+0	E
Y0948	YR109c	337-193982-194572		0.25 1091 (1596)	hypothetical protein, similar to YKR042c (PIR: S38114), TMM 1+0	E
Y0955	YR107a	279-195985-196823	ASF1	0.14	ASF1 protein (PIR: S30766), TMM 1+1	A
Y0960		197011-197083			rDNA ²⁰	
Y0963		197085-197242			E element	
Y0970		197245-197613			cdc 1, LTR of Ty4	
Y0971		444-197615-198054	Ty44_II	0.17	Ty4A_II protein	
Y0980		1803-197615-203022	Ty4B_II	0.15	Ty4B_II protein	
Y0983		203026-203484			cdc 1, LTR of Ty4	
Y0986		203503-203814			E element	
Y0993		203815-204092			E element	
Y0996		204231-204502			rDNA ²⁰	
Y0997	YR102a	714-205081-207142		0.12 329 (1303)	probable G-protein, β -transducin type (PIR: S48886)	D
Y0998	YR101a	530-205573-208222		0.09 1754 (1527)	probable component of the TCP-1 ring complex, similar to mouse CCT7 (PIR: S43056)	C
Y0999	YR109c	551-206621-211273	GDF3	0.30 274 (2405)	GATA zinc finger protein 3 (GB: X66331)	B
Y0998	YR109c	1769-210499-212005		0.17	hypothetical protein, TMM 5+1	E
Y0999	YR109c	385-211484-211972		0.17	hypothetical protein, TMM 8+1	E
Y0999	YR105c	367-213552-213712		0.13	hypothetical protein	F
Y0999	YR109a	645-221896-223020	SMO2	0.15	probable protein kinase SMO2 (PIR: S20136), TMM 1+0	A
Y0999	YR105a	580-224751-226439		0.10 386 (2134)	hypothetical protein, similar to YKR027a (PIR: S38101), TMM 1+0	E
Y0999	YR104a	649-227023-227469		0.09	hypothetical protein, 7	F

Table B. Continued

Nomenclature	Size (aa)	Coordinates	Locus	CAI	PaaA score	Description (nature of element, function or similarity of product)	Comment
Working Official							
R025		228122-228297	<i>SNR57</i>			<i>snR 37</i> small nuclear RNA	
R024	YJL005	618-228724	<i>SNR77</i>	0.12	253 (2086)	probable heat dependent regulatory protein, similar to S0414	D
R026	YJL024	819-230997	<i>MEF2</i>	0.13		translation elongation factor G homolog, MEF2, mitochondrial (PIR: S47548), TMM 1+1	A
R029		233635-233707				<i>snR4TM</i>	
R032	YJL011	678-234019	<i>CS01</i>	0.14		glutamate-cysteine ligase (PIR: S28648), TMM 2+1	A
R034	YJL006	607-234959	<i>CS079</i>	0.11		hypothetical protein	F
R038	YJL009	746-235110	<i>CS03</i>	0.12		CS03 protein (GB: U15905)	A
R040	YJL004	1058-241778	<i>CS051</i>	0.15	1625 (4985)	hypothetical protein, similar to YKR028a (GB: X05021)	F
R062	YJL007	217-245287	<i>CS037</i>	0.18		hypothetical protein, TMM 6+0	E
R064	YJL004	224-245997	<i>CS068</i>	0.15		hypothetical protein, TMM 2+0	E
R066	YJL004	1478-246950	<i>BCK1</i>	0.12		protein kinase BCK1 (PIR: S20137)	A
R069	YJL004	873-250519	<i>CS037</i>	0.15	264 (4296)	probable transport protein, similar to PIR: A42111, TMM 13+0	E
R081	YJL003	691-254435	<i>TRK1</i>	0.12		<i>TRK1</i> , outwardly rectifying potassium channel protein, TMM 20+0 F	
R083	YJL025	1174-257118	<i>RAD59</i>	0.13		telomere RAD59 (PIR: S48586)	A
R085	YJL001	498-260778	<i>CS0271</i>	0.13		hypothetical protein, TMM 5+1	E
R088	YJL000	784-262455	<i>CS0246</i>	0.14		hypothetical protein	F
R092	YJL009	829-263621	<i>S09</i>	0.14		<i>S09</i> protein, probable regulatory protein (GB: U17843), TMM 2+1	A
R094	YJL009	480-264388	<i>S0907</i>	0.16		ornithine carbonyltransferase (PIR: S00056), TMM 1+1	A
R097	YJL001	827-264708	<i>PRX1</i>	0.16		<i>snR4</i> ligase (PIR: A29617), TMM 1+0	A
R099	YJL006	122-272136	<i>CS0241</i>	0.11		hypothetical protein, TMM 1+0	E
R099	YJL001	623-272522	<i>CS0290</i>	0.16		hypothetical protein	F
R094	YJL004	1048-274980	<i>TRP97</i>	0.13	1595 (4883)	hypothetical protein, similar to YKR021W (PIR: S38995)	F
R102	YJL003	604-278736	<i>CS0247</i>	0.09	596 (2822)	hypothetical protein, similar to YKR019 (PIR: S3888)	F
R107	YJL001	731-280880	<i>CS032</i>	0.17	2652 (5986)	hypothetical protein, similar to YKR016 (PIR: S3887), TMM 1+1	E
R102	YJL001	489-283500	<i>ACT3</i>	0.13		actin-related protein (PIR: S47808)	A
R107	YJL000	1221-285236	<i>SCP190</i>	0.15		<i>SCP190</i> protein, histone-like protein (PIR: S37492)	A
R102	YJL000	298-289575	<i>CS0489</i>	0.30	679 (1288)	hypothetical protein, similar to YKR017W (PIR: S38862), TMM 3+0	C
R107	YJL000	881-291634	<i>CS0676</i>	0.15	597 (3522)	hypothetical protein, similar to YKR017W (PIR: S38862), TMM 2+0	D
R103	YJL001	131-294364	<i>CS0276</i>	0.08		hypothetical protein, TMM 1+1, 7	E
R108	YJL000	1189-296980	<i>CS0506</i>	0.15	505 (4986)	potato protein-binding protein, similar to YKR016 (PIR: S25414)	D
R104	YJL001	138-298338	<i>CS0571</i>	0.11		hypothetical protein, TMM 1+0	E
R104	YJL004	1230-298835	<i>CS0544</i>	0.18	605 (5561)	probable protein nucleotide-binding protein, similar to SMC1 (PIR: S43966), TMM 1+0	D
R100	YJL003	602-302735	<i>CS0480</i>	0.14		hypothetical protein, TMM 1+1	E
R106	YJL001	213-304919	<i>CS0557</i>	0.12		hypothetical protein, TMM 1+0	E
R100	YJL001	374-305827	<i>CS0548</i>	0.12	514 (2803)	similar to acetylglutamate synthase (GB: U35884), TMM 1+1	D
R100	YJL000	888-307869	<i>CS0532</i>	0.18	461 (4614)	hypothetical protein, similar to YKR264a (PIR: S47120), TMM 1+1	E
R108	YJL000	594-310820	<i>CS0481</i>	0.17		hypothetical protein	F
R102	YJL000	299-312114	<i>CS0610</i>	0.20	525 (1572)	similar to human annexin D (SW: P10766)	D
R107	YJL000	116-313779	<i>CS0426</i>	0.12		hypothetical protein, TMM 1+1	E
R111	YJL000	252-313812	<i>CS0467</i>	0.16		hypothetical protein	F
R115	YJL000	167-314732	<i>CS0525</i>	0.11		hypothetical protein	F
R120	YJL000	131-314870	<i>CS0262</i>	0.12		hypothetical protein, TMM 1+1	E
R125	YJL000	218-315437	<i>CS0470</i>	0.09		ribosomal protein L17, mitochondrial (PIR: S47128)	A
R132	YJL002	830-316979	<i>CS0468</i>	0.12		hypothetical protein, TMM 9+1	E
R133	YJL001	713-319711	<i>CS0449</i>	0.16		hypothetical protein	F
R138	YJL000	444-320981	<i>CS0402</i>	0.21	462 (2151)	probable amino acid transferase, similar to (PIR: S52796)	D
R139	YJL000	408-324839	<i>CS0482</i>	0.12		hypothetical protein, TMM 6+1	E
R141	YJL000	343-325940	<i>CS0568</i>	0.12	1119 (2463)	protein nucleotide binding protein, similar to YKR270 (PIR: S46111), TMM 1+0	C
R143	YJL001	663-327816	<i>CS0608</i>	0.18		hypothetical protein, TMM 1+1	E
R143	YJL000	880-330129	<i>CS0708</i>	0.16	436 (4257)	probable regulatory protein, similar to mouse Ku2 protein (PIR: S08549), leucine zipper D	
R148	YJL000	240-333832	<i>CS0786</i>	0.14		hypothetical protein	F
R150	YJL000	478-335960	<i>CS0883</i>	0.15		hypothetical protein	F
R152	YJL000	379-335995	<i>CS0729</i>	0.14		<i>PEP5</i> protein (PIR: S48882)	A
R154	YJL002	332-337966	<i>CS0861</i>	0.06		glyceraldehyde-3-phosphate dehydrogenase 3 (PIR: A0072), TMM 1+1	A
R158	YJL001	622-338482	<i>CS0947</i>	0.12		hypothetical protein, TMM 3+0	E
R158	YJL000	1973-342217	<i>CS0439</i>	0.20	971 (3214)	oral mRNA translation inhibitor SK12 (GB: U29641)	D
R162	YJL000	450-345668	<i>CS0707</i>	0.16		hypothetical protein	F
R164	YJL000	396-347145	<i>CS0532</i>	0.14	344 (1921)	hypothetical protein, similar to YKR273 (PIR: S46134)	F
R166	YJL001	842-349278	<i>CS0803</i>	0.12		hypothetical protein	F
R171	YJL000	451-351935	<i>CS0507</i>	0.12	302 (2257)	similar to Igmu protein ligase A E-cad (PIR: A54031)	D
R173		353839-354027				<i>snR4TM</i> (small intron)	
R177		354235-354555				<i>snR 5</i>	
R179		354539-354870				<i>snR 6</i>	
R180		355069-355140				<i>snR4TM</i>	
R180		355131-355222				<i>snR4TM</i>	
R184	YJL000	634-357119	<i>CS0620</i>	0.16	2721 (3048)	similar to succinate dehydrogenase flavoprotein (PIR: S34793)	B
R202	YJL000	458-357998	<i>CS0571</i>	0.16		GTPase-activating protein GTP (PIR: S3066), TMM 1+0	A

Table II. Continued

Nomenclature	Size (aa)	Coordinates	Locus	CAI	Fasta score	Description (nature of element, function or similarity of product)/Comment	
Working Official							
YJ204	YJL041w	257	359825-360590	0.09		hypothetical protein	F
YJ206	YJL042w	158	360644-361137	0.15		nucleobin-associated protein (GB: XNMF2)	B
YJ207	YJL041w	825	365479-366065	0.16		nucleobin-like protein NSP1 (PR: S14055) (clone from M5480 to M5597)	B
YJ216	YJL039w	5683	366446-373494	0.15		hypothetical protein, TMM 4+1	E
YJ221			374139-374190			ORF TM	
YJ226			374204-374272			ORF TM	
YJ230			374579-374630			orf 2	
YJ232	YJL036w	219	374815-375469	0.10	405 (1049)	similar to YJ234, TMM 3+0	E
YJ234	YJL037w	224	376297-377028	0.11	405 (1049)	similar to YJ232, TMM 2+1	E
YJ240			378055-378128			ORF TM	
YJ244	YJL036w	423	378720-379798	0.15		hypothetical protein	F
YJ246	YJL035w	250	379847-380696	0.12		hypothetical protein	F
YJ249	YJL034w	682	381027-383067	0.44		nuclear fusion protein KAR2 precursor (PR: A32366), TMM 1+1	A
YJ250	YJL033w	779	383532-385861	0.20	530 (1629)	similar to <i>Levini</i> SemB RNA helicase (SW: F23207)	D
YJ252	YJL032w	194	386645-386764	0.15		hypothetical protein	F
YJ254	YJL031w	290	386966-388035	0.15		propylgeranyl transferase α chain (PR: S48302)	A
YJ256	YJL030w	196	387152-387939	0.12		MAO2 protein (PR: S48302)	A
YJ259	YJL029w	822	388955-390548	0.15	517 (1644)	similar to <i>C. elegans</i> T09G5.8 protein (PR: S41008)	F
YJ263			390736-390800			ORF TM	
YJ267	YJL028w	111	391006-391138	0.07		hypothetical protein, TMM 2+0, 1	E
YJ269	YJL027w	139	391531-391644	0.08		hypothetical protein, 1	F
YJ271	YJL026w	309	392099-393295	0.50		ribonucleoside-diphosphate reductase small chain (PR: A26936), TMM 1+1	A
YJ273	YJL025w	514	393662-395203	0.13		RRNT protein (PR: S30785)	A
YJ274	YJL024w	194	395623-396287	0.14	229 (708)	related to mouse clathrin associated protein 19 (clone from 396399 to 396265) (PR: A80311)	D
YJ278			396421-396491			ORF TM	
YJ282	YJL023w	541	397055-398693	0.15		hypothetical protein	F
YJ284	YJL022w	192	397866-398189	0.10		hypothetical protein, TMM 1+1, 1	E
YJ286	YJL021w	585	398635-399729	0.13		hypothetical protein	F
YJ287	YJL020w	771	399709-402191	0.14	206 (368)	glutamic acid rich protein precursor (<i>Flaccodon</i> <i>flaccodon</i>) (PR: A54514)	D
YJ290	YJL019w	620	402388-404647	0.12		hypothetical protein, TMM 1+0	E
YJ293	YJL018w	104	404021-404632	0.16		hypothetical protein	F
YJ295	YJL017w	325	405276-406252	0.15		hypothetical protein	F
YJ296	YJL016w	171	406447-406699	0.16		hypothetical protein	F
YJ331	YJL015w	124	408834-407205	0.12		hypothetical protein	F
YJ336	YJL014w	534	407246-408647	0.25		chaperonin of the TCP-1 ring complex, TMM 1+1, similar to mouse CCT3 (PR: S49062)	B
YJ341	YJL013w	515	409186-410728	0.15	475 (2456)	similar to protein kinase BLUB1 (Yeast chr 1) (GB: LMF2027)	D
YJ345	YJL012w	648	411143-413086	0.25		hypothetical protein	F
YJ349	YJL011w	161	413975-414407	0.12		hypothetical protein	F
YJ352			414655-414725			ORF TM	
YJ353			415616-415724			ORF TM (small insert)	
YJ357	YJL010w	666	417232-419249	0.17		hypothetical protein	F
YJ369	YJL009w	108	419542-419665	0.16		hypothetical protein, TMM 1+1	E
YJ374	YJL008w	568	419647-421350	0.20	1219 (2622)	probable chaperonin of the TCP-1 ring complex, similar to mouse CCT8 (PR: S32887)	C
YJ379	YJL007w	164	422388-422699	0.15		hypothetical protein, TMM 1+0	E
YJ385			422624-422696			ORF TM	
YJ390	YJL006w	323	422828-423796	0.11		hypothetical protein, TMM 1+0	E
YJ395			426119-426202			ORF TM	
YJ401	YJL005w	2026	426846-430021	0.12		adenylate cyclase (PR: A24776)	A
YJ402	YJL004w	263	431279-431687	0.09		hypothetical protein, TMM 4+0	E
YJ403	YJL003w	118	432331-432664	0.10		hypothetical protein, TMM 1+0, 1	E
YJ404	YJL002w	476	432941-434038	0.16		α subunit, oligomannosyltransferase (GB: Z06710), TMM 2+0	A
YJ407	YJL001w	193	435032-435630	0.17		multicatalytic endopeptidase complex chain PRE3 (PR: S43669), TMM 1+0	A
			435996-436028			commonest	
			436022-436104			commonest	
			436105-436162			commonest	
YJ409	YJL000w	602	436489-438294	0.12	237 (2951)	similar to <i>C. elegans</i> , hypothetical protein (PR: S42372), TMM 10+1	E
YJ411	YJL000w	593	438531-440329	0.17		hypothetical protein	F
YJ415	YJL000w	539	440885-442399	0.15		hypothetical protein	F
YJ418	YJL000w	650	442598-444547	0.13		α -agglutinin (PR: S12835), TMM 2+0	A
YJ422	YJL000w	445	445689-447708	0.19		clathrin-associated protein complex β chain homolog (PR: S12554), TMM 1+1	A
YJ427	YJL000w	487	448008-450528	0.16		hypothetical protein	F
YJ429	YJL000w	364	450706-451617	0.17		translational initiation factor eIF-2 α chain (PR: A32108)	A
YJ431	YJL000w	338	452116-453129	0.14		hypothetical protein	F
YJ433	YJL000w	332	453372-454367	0.10		glyoxaldehyde 3-phosphate dehydrogenase (PR: S40911)	A
YJ436	YJL000w	511	455925-457497	0.29		sulfate adenylyltransferase (PR: S09906)	A
YJ438	YJL000w	261	458139-459152	0.14		hypothetical protein	F
YJ440	YJL000w	267	459464-460054	0.12		hypothetical protein, TMM 1+0	E

Table II. Continued

Nomenclature	Size (aa)	Coordinates	Locus	CAI	Fasta score	Description (nature of element, function or similarity of product)/Comment		
Working Official								
J1444	YB0115a	305	460365	461271	0.11	hypothetical protein, TMM 3+1	E	
J1446	YB0115a	798	461516	462109	0.22	hypothetical protein	F	
J1448	YB0115a	518	462408	463103	0.13	1980 (2637)	similar to SNG1 gene (yeast chr 7) (GB: X14926), TMM 3+1	C
J1450	YB0115a	585	464041	464895	0.26	dihydroxy-acid dehydratase (PDB: 5AT9a)	A	
J1452	YB0115a	190	466211	466780	0.12	ESS1 protein (PDB: 507807)	A	
J1454	YB0115a	126	466475	466632	0.08	hypothetical protein, TMM 1+1, 1	E	
J1456	YB0115a	349	466822	467368	0.11	222 (17786)	similar to E. coli <i>usp-C</i> gene	D
J1458	YB0115a	119	467688	468007	0.11	hypothetical protein, TMM 1+1	E	
J1462	YB0115a	292	468310	469266	0.11	MEK2, meiotic recombination protein MEK2 (clone from M987) or M985 (PDB: 1A627)	A	
J1464	YB0115a	128	469416	469787	0.13	hypothetical protein	F	
J1470	YB0115a	153	469896	469992	0.09	hypothetical transport protein, TMM 2+1, 1	E	
J1503	YB0125a	244	469820	470651	0.12	hypothetical protein	F	
J1550	YB0125a	171	470826	471358	0.17	313 (922)	similar to human 3-hydroxybutyrate 3,4-dehydrogenase (PDB: 1S470)	D
J1553			472150	472487		3, LTR of Ty1		
J1555		440	472447	473786	0.14	1990 (2005)	TyA protein	
J1560		1741	472447	477702	0.15	8241 (8276)	TyB protein	
J1563			477736	478071		3, LTR of Ty1		
J1565		440	478031	479330	0.15	1991 (1997)	TyA protein	
J1570		1741	478031	483296	0.14	8251 (8277)	TyB protein	
J1573			483322	483659		3, LTR of Ty1		
J1575	YB0135a	745	483649	485683	0.11	443 (1553)	hypothetical protein, similar to 30450	F
J1580	YB0135a	1408	486276	490499	0.15	3171 (3683)	hypothetical protein, similar to Y12.022a (PDB: 5Z036), TMM 6+1	E
J1583	YB0135a	393	490768	491946	0.19	468 (1962)	hypothetical protein, similar to LR167.1a (PDB: 548567)	F
J1586	YB0135a	1357	492066	496138	0.14	3055 (3711)	hypothetical protein	F
J1604	YB0135a	108	496370	496641	0.12	PEY191 protein (PDB: 5Z924)	A	
J1606	YB0135a	1085	497642	500296	0.13	probable helicase RAN100 (SW: P40372), TMM 1+1	A	
J1608	YB0135a	892	500403	503076	0.11	hypothetical protein, TMM 1+1	E	
J1610	YB0135a	127	502789	503189	0.11	hypothetical protein	F	
J1612	YB0135a	120	503400	503739	0.08	hypothetical protein, TMM 2+0, 1	E	
J1614	YB0135a	1121	503625	506885	0.13	hypothetical protein, TMM 2+1	E	
J1616	YB0135a	779	507410	509769	0.14	788 (3954)	similar to mouse chloride channel protein (GB: D17521), TMM 7+1	D
J1622	YB0135a	1174	509920	513450	0.14	hypothetical protein, TMM 2+1	E	
J1624	YB0135a	744	513742	515973	0.13	hypothetical protein, TMM 1+0	E	
J1626	YB0135a	750	516150	517200	0.14	hypothetical protein	F	
J1631			517500	517571		rRNA ^{23S}		
J1634			517605	517786		3, rRNA ^{23S}		
J1637	YB0135a	140	518451	518872	0.15	hypothetical protein, TMM 4+0	F	
J1639	YB0135a	454	519326	521289	0.52	heat shock protein 70-related protein SBC1 precursor, mitochondrial (PDB: A3249)	A	
J1641	YB0135a	604	521735	523546	0.11	hypothetical protein, TMM 1+1	E	
J1647			523699	523780		rRNA ^{23S}		
J1651	YB0135a	157	524798	525068	0.20	translation initiation factor eIF-3A.2 (PDB: B40259)	A	
J1653	YB0135a	109	526022	526348	0.37	cytochrome c isoform 1	A	
J1655	YB0135a	530	526574	528365	0.13	UTR1 protein (PDB: 546599), TMM 1+1	A	
J1657	YB0135a	235	528384	529088	0.10	UTR1 protein (PDB: 546599)	A	
J1659	YB0135a	504	529548	531050	0.17	ORF1 protein precursor (PDB: 546791), TMM 1+0	A	
J1661			531202	531361		3, rRNA ^{23S}		
J1663			531513	531585		rRNA ^{23S}		
J1665	YB0135a	565	531749	533445	0.14	RAD7 protein (PDB: A2523a)	A	
J1667	YB0135a	576	533764	535435	0.15	hypothetical protein	F	
J1669	YB0135a	497	535743	537210	0.17	325 (2446)	hypothetical protein, similar to YNP027.05a (GB: Z17916), TMM 4+0	E
J1670			538242	538313		rRNA ^{23S}		
J1705	YB0135a	164	538499	538950	0.13	HIT1 protein (PDB: 538866)	A	
J1706a			540451	540783		5, rRNA ^{23S}		
J1706b			540786	541114		5, rRNA ^{23S}		
J1707			541195	541266		rRNA ^{23S}		
J1710	YB0135a	236	541482	542289	0.10	hypothetical protein	F	
J1713			542643	542730		rRNA ^{23S} (small subunit)		
J1715	YB0135a	216	543749	544796	0.15	GMP kinase (PDB: A06683)	A	
J1720	YB0135a	147	544422	544962	0.08	clathrin-associated protein 17 (PDB: C40535)	A	
J1725	YB0135a	818	545474	547927	0.16	1251 (3786)	similar to vertebrate-specific protein kinase (PDB: 538055), TMM 1+0	D
J1730	YB0135a	351	548426	549498	0.14	continuous-binding protein CP1 (PDB: A36310)	A	
J1736	YB0135a	935	550198	553000	0.13	hypothetical protein, TMM 1+1	E	
J1742	YB0135a	417	553186	554536	0.12	amino-terminal nucleic acid (PDB: 547936)	A	
J1747	YB0135a	125	554882	555256	0.20	DNA-directed RNA polymerase I chain A12.2 (PDB: A48107), TMM 1+0	A	
J1752	YB0135a	562	555601	557386	0.22	1704 (2677)	probable chaperonin of the DCP-1 ring complex, similar to mouse CCT3 (PDB: 543981), TMM 1+0	C
J1760	YB0135a	449	557199	558845	0.20	1499 (2151)	similar to actin-like protein Act 2 (baker's yeast) (PDB: A11786), TMM 1+1	C
J1800	YB0135a	2470	559103	564512	0.14	YOR1 protein (PDB: 547940), TMM 1+1	A	
J1805	YB0135a	141	566709	567131	0.14	hypothetical protein	F	

Table II. Continued

Nomenclature	Size (aa)	Coordinates (aa)	Locus	CAI	ProtA score	Description (nature of element, function or similarity of product)/Comment	
Working	Official						
Y0808	Y0808a	335	561336-561388	RFC2	0.18	replication factor C chain RFC2 (PR: 54551)	A
Y0811	Y0809c	197	561406-561606		0.20	hypothetical protein	F
Y0814	Y0807a	325	561511-570265		0.40	hypothetical protein	F
Y0818	Y08071a	122	570002-570437		0.08	hypothetical protein, 7	F
Y0821	Y0807c	309	570657-571411		0.17	827 (1844)	F
Y0824	Y0807b	206	572005-572622	PKM2	0.17	methylester fatty-acyl-phospholipid synthase (PR: R20441), TMM 1+1	A
Y0827	Y08074a	218	572782-573135		0.15	hypothetical protein	F
Y0830	Y08075a	396	573668-574635		0.18	289 (202)	F
Y0833	Y08076c	415	575044-576208	CDC31	0.17	cell division control protein CDC31 (PR: 54910)	D
Y0837	Y08077c	311	576625-577677	MRP7	0.36	phosphate transport protein, mitochondrial (PR: 512116), TMM 1+1	A
Y0840	Y08078a	473	578247-579969		0.13	914 (225)	F
Y0843	Y08080a	109	579952-580025		0.10	hypothetical protein, common from Y08055 to Y08079a, TMM 1+0	D
Y0847	Y08081c	394	580172-581303		0.14	hypothetical protein	F
Y0854	Y08082c	103	581604-581642		0.15	hypothetical protein	F
Y0857	Y08083c	309	582298-583224		0.10	hypothetical protein	F
Y0860	Y08084a	423	583426-584688		0.19	hypothetical protein	F
Y0863	Y08085c	105	584816-585024		0.14	hypothetical protein, TMM 2+0	E
Y0866	Y08086a	100	585755-586062	STE18	0.18	STE18 protein (PR: R3002)	A
Y0870	Y08087a	108	586087-586434		0.16	hypothetical protein, TMM 2+0, 7	E
Y0875	Y08088c	292	586185-587080		0.17	hypothetical protein	F
Y0880	Y08089a	354	587497-590296		0.13	hypothetical protein	F
Y0883	Y08090c	103	590362-590614	GRR1	0.12	GRR1 protein (PR: A40226), TMM 1+1	A
Y0890	Y08091c	109	591751-591825		0.15	793 (4842)	F
Y0891	Y08091b	289	592187-593035		0.15	ATP/GTP binding site motif A	E
Y0895	Y08092a	1320	598009-602708		0.15	hypothetical protein	F
Y0911	Y08093c	527	602916-603996	POP1	0.12	component of pre-mRNA polyadenylation factor	B
Y0916	Y08094c	560	604265-605344	IME1	0.18	meiosis-inducing protein IME1 (PR: 57117)	A
Y0923	Y08095a	522	609466-610431	ACR1	0.20	ACR1 protein (PR: 543280), TMM 2+1	A
Y0929	Y08096a	287	610888-611733		0.22	431 (149)	D
Y0931	Y08097a	172	612096-612623		0.13	hypothetical protein	F
Y0936	Y08098c	655	612882-614846		0.15	hypothetical protein	F
Y0941	Y08099a	236	615266-615933	YU81	0.11	ubiquitin carboxyl-terminal hydrolase YU81 (PR: 51132), TMM 1+0	A
Y0946	Y08100c	527	616644-617624		0.18	hypothetical protein	F
Y0950			617600-617700			RNA ^{5S} -small protein	
Y0952	Y08101a	286	617924-618721		0.11	hypothetical protein	F
Y0957	Y08102c	202	618850-619295		0.13	hypothetical protein	F
Y0962	Y08103a	564	620440-621335	URX4	0.16	CTP synthase URX4 (PR: 542780), TMM 2+0	A
Y0968	Y08104c	154	622242-622703	SKO1	0.38	spermidine deiminase (Ca ²⁺) (PR: A3617)	A
Y0973	Y08105a	540	623270-624289		0.15	hypothetical protein	F
Y0978	Y08106a	729	624927-626764		0.18	hypothetical protein	F
Y0983	Y08107a	526	627050-628037		0.13	hypothetical protein, TMM 12+1	E
Y0988	Y08108a	123	628400-628771		0.14	hypothetical protein	F
Y0992	Y08109c	1119	629270-632672	CPB2	0.24	large subunit of arginine specific carbamoylphosphate synthase (PR: A61099)	A
Y0997	Y08110a	648	633306-633589	CPB1	0.16	small subunit of arginine specific carbamoylphosphate synthase (PR: R13478)	A
Y0999	Y08111c	263	635549-636397		0.12	hypothetical protein	F
Y0999	Y08112a	203	636721-637323		0.09	hypothetical protein	F
Y0999	Y08113c	247	637926-638666		0.10	204 (118)	F
Y0999	Y08114a	136	638350-638739		0.11	hypothetical protein, TMM 1+0	E
Y0999	Y08115a	189	639633-640139		0.10	hypothetical protein	F
Y0999	Y08116a	278	640516-641392		0.14	hypothetical protein, TMM 2+1	E
Y0999	Y08117a	453	641898-643056		0.25	hypothetical protein, TMM 5+1	E
Y0999	Y08118c	207	643084-643792		0.19	hypothetical protein, TMM 3+1	E
Y0999	Y08119c	728	643998-646181		0.15	778 (3628)	F
Y0999	Y08120a	116	646617-647164		0.07	similar to human retinoblastoma binding protein 2 (PR: 506411)	D
Y0999	Y08121a	311	647298-648036	APP2	0.42	ATP synthase β chain precursor (PR: 52729)	A
Y0999	Y08122a	497	649467-650667		0.15	hypothetical protein	F
Y0999	Y08123a	125	651592-651296	RP57	0.75	ribosomal protein S5	A
Y0999	Y08124c	448	652986-653924		0.14	hypothetical protein, TMM 0+1	E
Y0999	Y08125c	408	654431-655654		0.17	283 (177)	F
Y0999	Y08126c	811	655948-658168		0.13	321 (198)	F
Y0999	Y08127c	1360	658621-662750	ZMS1	0.12	ZMS1 protein (PR: 543714), TMM 4+1	A
Y0999	Y08128a	119	662817-662968		0.09	hypothetical protein, 7	F
Y0999			663440-663673	SNR7		SNR 7 small nuclear RNA	
Y0999	Y08129c	139	663684-664710		0.11	hypothetical protein, TMM 1+0	E
Y0999	Y08130c	679	664912-666928		0.13	1758 (107)	F
Y1010	Y08131a	349	667035-668990	MNS1	0.14	similar to human prostate-specific membrane antigen (PR: Q09006), TMM 1+0	D
Y1017	Y08132a	348	669213-672158		0.15	hypothetical protein, TMM 2+1	E
Y1018	Y08133a	208	672682-673308		0.28	hypothetical protein	F
Y1020	Y08134c	307	673423-673947		0.13	343 (1882)	F

Table II. Continued

Nomenclature	Size (aa)	Coordinates (aa)	Locus	CAI	FastA score	Description (nature of element, function or similarity of product)/Comment	
Working Official							
J2121	YJR115c	239	675155-676489	0.12		hypothetical protein	F
J2124	YJR116c	421	677115-678297	0.10		hypothetical protein	F
J2126	YJR117c	1447	679671-682976	0.25	1054 (8897)	similar to formyltransferase reductase (SW: P30098)	D
J2129	YJR118c	1584	684258-689009	0.14		hypothetical protein	F
J2132	YJR119c	109	689139-690215	ACM60	0.47	homocysteine dehydrogenase (PIR: S11107), TMM 1+1	A
J2161	YJR140c	1648	690444-695387	0.11		hypothetical protein, TMM 1+1	E
J2166	YJR141c	347	695597-696637	0.13		hypothetical protein, TMM 1+1	E
J2171	YJR142c	342	696832-697857	0.15		hypothetical protein	F
J2176	YJR143c	762	698020-700005	PM5V	0.22	PM14 protein (PIR: S51284), TMM 8+1	A
J2181	YJR144c	269	700575-701379	MGM107	0.16	MGM104 protein (PIR: S26449)	A
J2186	YJR145c	261	701721-702799	RP374	0.69	ribosomal protein S4c; 10 clones from 702490 to 702740 (PIR: S28054)	A
J2200	YJR146c	117	703336-703926	0.07		hypothetical protein, 7	F
J2204	YJR147c	358	703887-704960	0.12	277 (1762)	similar to heat shock transcription factor 8 (PIR: S25481)	D
J2209	YJR148c	376	705435-706062	0.19	1584 (1988)	similar to TW11 yeast protein (PIR: S49999)	C
J2213	YJR149c	404	706851-708062	0.18	462 (1937)	similar to 2-nitropropane dioxygenase (PIR: S58911)	D
J2217	YJR150c	298	708205-709298	0.30		hypothetical protein, TMM 2+0	E
J2223	YJR151c	1161	711949-715431	0.23	814 (4582)	similar to human myosin (PIR: A49963), TMM 2+0	D
J2230	YJR152c	543	719337-720985	DMJ7	0.16	alkaline phosphatase (PIR: A26671), TMM 8+1	A
J2233	YJR153c	361	722908-723588	0.17	907 (1643)	similar to polygalacturonase (PIR: S26771), TMM 1+0	C
J2240	YJR154c	346	725415-726012	0.13		hypothetical protein	F
J2245	YJR155c	288	727636-727959	0.15	1334 (1479)	similar to yeast acylalcohol dehydrogenase (PIR: S11315)	B
J2250	YJR156c	340	728268-729287	0.21	1784 (1790)	similar to thiamine-repressed orn-1 protein (PIR: S48548), TMM 1+0	B
J2253	YJR157c	120	730206-730585	0.13		hypothetical protein, TMM 1+0	F
J2260	YJR158c	367	732131-733051	0.16	1893 (3036)	similar to hexose transport protein HXT7 (PIR: S43186), TMM 9+1	C
J2265	YJR159c	317	733735-734685	ROMV	0.22	oxalate dehydrogenase (GI: L13039)	B
J2400	YJR186c	602	737302-739587	0.13	2583 (8048)	similar to sugar transport protein (SW: P30356), TMM 7+1	C
J2410	YJR184c	383	742542-743690	0.14	1643 (2635)	similar to YBR6 (SW: P30363), TMM 3+1	E
			744595-745052			core X element	
			745053-745558			STR-D, C, B and A elements	
J2420	YJR192c	116	744685-744952	0.14	422 (806)	similar to YKWS (SW: P36576) right intronic sequence	F

Last column: status of the protein deduced from each putative gene. The categories A (fully known) to F (unknown) are defined in the text. The self FastA score of the predicted protein is in parentheses. An accession number in one of the public databases (PIR, Swiss-Prot (SW), GenBank (GB) and EMBL) is indicated. Abbreviations: TMM, transmembrane motif, integral + peripheral; ? , questionable gene. ORF YJL093c is categorized as F, as it was discovered and sequenced during the systematic sequencing of chromosome X and found to correspond to no known gene. It was subsequently biologically characterized as a potassium channel (Kerchum et al., 1995).

novel putative yeast genes whose function will have to be determined experimentally. However, 57 of these (another 15% of total) encode proteins that show significant similarity to a protein of known function from yeast or other organisms, thus providing some indication as to their function. The 204 (54%) remaining ORFs exhibit no significant similarity to known sequences (FastA score < 200). Motif searches have shown that 91 of the latter have some particular protein signature, mostly a structure suggestive of transmembrane domains (Table II).

An approximately equal number of ORFs is observed on each DNA strand. The mean ORF size is 482 codons (1446 bp), the longest (YJR066w) reaching 2470 codons. The mean size of inter-ORF regions, disregarding one in each pair of overlapping ORFs, is 602 bp for terminator-promoter combinations (WW and CC in Figure 3). For divergent promoters (DP) and convergent terminators (CT), the mean size is 725 bp and 311 bp, respectively. This striking difference in inter-ORF size between divergent promoters versus convergent terminators may be indicative of more important sequence requirements in promoter regions for the regulation of gene expression. An exception is the contiguity of the two ORFs YJL108c and YJL107c. The TGA stop codon of the latter overlaps the ATG of the former, so that both codons share TG. This peculiarity was carefully checked by oligo-primed sequencing in

either direction on cosmid DNA. The two ORFs in their integrity are translated from a single transcript of ~3 kb (Rasmussen, 1995).

Environment of ATG and stop codons

Compilation of a large number of sequence data surrounding the initiation codon AUG has revealed that these sequences are not random and that higher eukaryotes have in common the consensus sequence GCC(A/G)CCATGG (Kozak, 1987). In the case of the budding yeast, another consensus (A/Y)A(A/Y)A(A/Y)AATGGTCT has been proposed (Hinnebusch and Liebman, 1991).

We examined the 318 chromosome X ORFs longer than 150 codons, in all probability corresponding to real genes, to test this consensus. Table III shows the frequency of the different nucleotides, as determined by tabulating positions -8 to +7 relative to ATG. A χ^2 test was performed at each position to test the non-randomness of this distribution, taking into account the G+C content of the chromosome. At all positions except -5 the distribution was found to be non-random. As these calculations are based on all the ORFs of a chromosome, regardless of their expression level, rather than on a selected subset, the following consensus sequence might be more appropriate: AAANAAAAATGGCTG. The chances of a random distribution at each position is < 5%, or even 1%

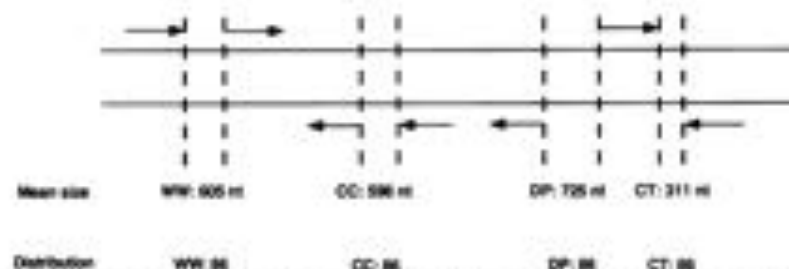


Fig. 3. Mean size and distribution of inter-ORF regions of chromosome X. WW: promoter/terminator combination on Watson strand; CC: promoter/terminator combination on Crick strand; DP: divergent promoters; CT: convergent terminators. The numbers indicate on top line the mean size, on bottom line the distribution of each configuration.

Table III. Initiation and stop codon environment

ATG environment													
	-8	-7	-6	-5	-4	-3	-2	-1	ATG	+4	+5	+6	+7
A	0.396	0.393	0.368	0.349	0.399	0.569	0.405	0.456	ATG	0.318	0.283	0.324	0.327
G	0.164	0.160	0.211	0.335	0.148	0.195	0.119	0.145	ATG	0.296	0.129	0.151	0.299
C	0.173	0.192	0.176	0.220	0.189	0.113	0.252	0.173	ATG	0.132	0.362	0.182	0.129
T	0.267	0.255	0.245	0.296	0.264	0.123	0.226	0.223	ATG	0.254	0.343	0.343	0.242
χ^2	7.978	9.616	10.015	7.370	10.060	104.811	30.264	27.741	ATG	20.165	61.227	8.750	22.695
TAG stop codon environment													
	-5	-4	-3	-2	-1	TAG	+4	+5	+6	+7	+8	+9	
A	0.380	0.268	0.310	0.394	0.296	TAG	0.408	0.282	0.380	0.437	0.366	0.282	
G	0.127	0.183	0.253	0.211	0.211	TAG	0.231	0.127	0.295	0.211	0.197	0.141	
C	0.183	0.197	0.169	0.083	0.113	TAG	0.113	0.197	0.183	0.056	0.169	0.259	
T	0.310	0.352	0.268	0.310	0.380	TAG	0.268	0.394	0.197	0.296	0.268	0.338	
χ^2	2.975	2.127	1.173	5.999	5.024	TAG	4.336	2.651	5.580	4.178	1.250	2.522	
TAA stop codon environment													
	-5	-4	-3	-2	-1	TAA	+4	+5	+6	+7	+8	+9	
A	0.368	0.296	0.387	0.452	0.361	TAA	0.297	0.316	0.368	0.355	0.297	0.393	
G	0.161	0.226	0.232	0.097	0.142	TAA	0.187	0.136	0.174	0.122	0.161	0.142	
C	0.200	0.239	0.129	0.155	0.181	TAA	0.129	0.200	0.148	0.168	0.271	0.155	
T	0.271	0.239	0.252	0.296	0.316	TAA	0.387	0.348	0.310	0.355	0.271	0.310	
χ^2	2.358	3.484	6.237	17.687	4.314	TAA	4.559	2.173	1.980	3.310	9.646	3.552	
TGA stop codon environment													
	-5	-4	-3	-2	-1	TGA	+4	+5	+6	+7	+8	+9	
A	0.348	0.304	0.402	0.424	0.261	TGA	0.347	0.315	0.304	0.390	0.315	0.272	
G	0.174	0.239	0.239	0.087	0.163	TGA	0.185	0.196	0.283	0.196	0.174	0.206	
C	0.185	0.120	0.152	0.196	0.163	TGA	0.109	0.109	0.163	0.196	0.152	0.185	
T	0.293	0.337	0.207	0.293	0.413	TGA	0.359	0.380	0.250	0.217	0.359	0.337	
χ^2	0.626	4.244	4.900	9.008	7.990	TGA	2.966	3.641	7.964	4.773	0.720	1.494	

The position relative to start or stop codon is indicated at the top of the columns. The numbers in the columns give the relative frequency of each base at each position. χ^2 tests were performed with three degrees of freedom (threshold for an α risk of 5% is 7.815 and for an α risk of 1% is 11.345). Expected frequencies used in χ^2 tests are A = 0.32, T = 0.32, G = 0.17 and C = 0.17 in non-coding regions, A = 0.32, G = 0.20, C = 0.19 and T = 0.28 in coding regions. Tabulation performed on 318 ORFs >150 codons.

at positions -3, -2, -1, +4, +5 and +7. We then addressed the question of the possible existence of a consensus sequence in the environment of the stop codons. Not surprisingly, TAA is the more frequently used stop codon: 155 ORFs longer than 150 codons have it, while 92 have

TGA and 71 TAG. When the nucleotide environment between positions -5 and +9 (position +1 being defined by the T of the stop signal) was tabulated, we observed the frequencies reported in Table III. It appears that, in the case of TAA, there is a bias at position -2, which is

more frequently than expected occupied by A and less frequently by G, and at position +8, where C is increased. In the case of TAG, at position -2 the frequency of C is depressed, while this nucleotide is nearly always absent from position +7. Finally, in the case of TGA, the distribution deviates from randomness at three positions, -2, -1 and +6.

Small ORFs (<100 codons)

The choice of a minimal length of 99 sense codons between the first ATG and the stop signal, which dates back to 1979 (Galibert *et al.*, 1979), probably owes more to the widely used decimal numbering system than to proper insight into biological mechanisms. However, as mentioned above, this size is warranted in the case of yeast (Dujon *et al.*, 1994). In simulation experiments in which chromosome length and nucleotide composition was varied, the chances that ORFs longer than 150 codons will exist and still not correspond to a real gene are negligible. Conversely, the chances that ORFs in the range 100–149 codons will have no biological significance increase in proportion to decreasing size. However, a size of 100 codons is no impassable limit and obviously some ORFs smaller than 100 codons correspond to genes and, for that matter, quite a few proteins shorter than 99 amino acids may not be accounted for by post-translational processing. An example is provided by the small proteolipids PMP1 and PMP2 (40 and 43 amino acids), on chromosomes III and V, respectively (Navarre *et al.*, 1992; 1994). Analysis of the chromosome X sequence has revealed 344 small ORFs 50–98 sense codons in size. Comparison of the deduced proteins with database entries shows that one of these, J0526 (106425–106706), corresponds to the gene encoding subunit VIII of ubiquinol-cytochrome *c* reductase (Hemrika *et al.*, 1993). It is a 94-amino acid protein whose coding gene has been hitherto overlooked. Another instance is YKR057w, which encodes a ribosomal protein of 87 amino acids. Some small ORFs, such as J1567 (479710–479952), J1564 (477910–478074) and J15591 (474126–474368) have perfect or nearly perfect matches with Ty retrotransposon proteins of longer size. These small ORFs most probably result from frameshift mutations, a rather common occurrence in these retroposons. Finally, significant similarity is observed between some small ORFs located in the subtelomeric region, such as J0210 (9452–9852), and similar elements located on other chromosomes (K-B110 on chromosome XI or LA75 on chromosome IX). The other small ORFs, displaying no significant homology with database entries, cannot simply be discarded, since some probably correspond to real genes. Examples in point are J0523 (105893–106060), J1153 (337859–338143), J2123 (676661–676924) and J1425 (448166–448444), all with CAIs >0.2. Clearly, a screening programme taking into account parameters such as the ATG and stop codon environment and the CAI must be developed to approach the question of their existence as genes.

Sequence duplications

We have analysed the nucleotide sequence of chromosome X for the occurrence of sequences demonstrating high similarity to other genes of chromosome X (intrachromosomal duplications) and to genes in other yeast chromo-

somes (interchromosomal duplications), both at the nucleotide and the amino acid level (Table IV). Some of the duplicated ORFs have been functionally characterized. These results confirm earlier observations on chromosomes XI (Dujon *et al.*, 1994) and II (Feldmann *et al.*, 1994) of the high level of internal genetic redundancy in the yeast genome. Moreover, in addition to duplication of individual genes, duplication of syntenic segments has also occurred, syntenic in the present context of intraspecies duplications meaning that two or more genes situated closely on the same chromosome have their homologous loci also located close together, with the same respective orientation, on the other chromosome. As a rule, the physical distance and the nucleotide sequence between two ORFs on the same syntenic segment are not conserved. However, some degree of intergenic sequence conservation can be observed in a few cases, as exemplified in Figure 4.

tRNAs and transposons

Twelve tRNA genes are found on each strand (Figure 5), a density somewhat higher than that observed in the previously sequenced yeast chromosomes. The 24 tRNAs can transfer 13 amino acids in all and include four tRNA^{Pro}, all identical with the same GTC anticodon; four tRNA^{Met}, two identical with TCT, one with ACG and one with CCT, the last two with minor sequence differences. Of the three tRNA^{His}, two are identical while the third exhibits slight differences. The two tRNA^{Trp} have an identical sequence and include the same GTA anticodon.

Upon folding, all the predicted tRNAs fit in readily with the clover-leaf model, regarding stem length as well as loop size. All the canonical bases are observed in all cases but one. The exception is tRNA^{Met} at position 517571, which exhibits an A, instead of T as in the canonical GTWC sequence. Careful checking of the sequence has shown that this ATC sequence does not result from sequencing errors. However, a cloning artefact at some point in the construction of the cosmid library cannot be ruled out at this stage.

While the clover-leaf model is basically respected, 46 non-canonical or unpaired bases are observable in the stems of this two-dimensional configuration. Thirty-nine correspond to a GT base pairing, three to TT and CA and one to GG. An example of such tRNA folding is presented in Figure 6. These observations cannot be ascribed to sequencing or cloning incidents, since they have been observed by different investigators all working on different cosmids. Furthermore, the reality of such pairings has been established by direct RNA sequencing on mature tRNA and by mutagenesis experiments (Pütz *et al.*, 1993). However, it is also true that in the case of plant mitochondrial tRNAs, some (but not all) mismatched base pairs are so edited as to generate a Watson-Crick pair in the mature tRNA (Maréchal-Drouard *et al.*, 1993). While this phenomenon is not yet documented in nuclear yeast tRNA, the possibility of a similar editing process, whereby some of the 46 mispairings mentioned above would be converted into conventional Watson-Crick pairs, cannot be dismissed without additional sequence data or structural studies at the tRNA level. An alternative hypothesis is that some of the predicted tRNAs actually correspond to inactive pseudogenes.

Four of the tRNA genes encountered in chromosome

Table IV. Related genes from chromosome X

Gene/ORF on chromosome X	Related gene/ORF on other chromosome ^a	Functional description ^b	aa identity % ^c	nt identity % ^d
YIL223c	PAU1(5)	PAU1 protein	96.7 (1-120)/120	96.7 (1-360)/360
YIL219c	LGT7 hexose transport protein	LGT7	97.9 (1-567)/567	98.4 (883-1701)/1701
YIL200c	ACD1(12)	similar to acornin hydratase	55.3 (35-782)/782	50.8 (6-2278)/2267
YIL196c	YCR0153c (3)	probable transport protein	65.0 (39-879)/881	68.1 (684-2387)/2643
YIL196c	YCR0154c (3)	similar to sterol isomerase SUR4	58.4 (16-310)/310	60.3 (70-891)/890
YIL191c (CRY2)	CRY1 (3)	ribosomal protein S14eB	96.3 (5-138)/138	92.0 (8-414)/414
YIL190c (RPS24)	L8098c (12)	ribosomal protein S15a	99.2 (1-130)/130	89.1 (1-390)/390
YIL164c (SRA3)	TPK2 (11)	cAMP-dependent protein kinase	84.5 (69-397)/397	73.0 (255-886)/1091
YIL159c (YUR1)	KTR2 (10)	YUR1 protein	66.3 (37-424)/426	64.3 (269-1250)/1284
YIL138c (TIF2)	TIF1 (11)	translation initiation factor eIF-2	100 (1-385)/385	99.3 (1-1185)/1185
YIL133c (MRS5)	MRS4 (11)	mitochondrial splicing protein	76.2 (23-312)/314	70.5 (119-875)/942
YIL099c (CSD3)	YKR027c (11)	CSD3 protein	42.3 (1-844)/1058	37.3 (1759-2238)/2238
YIL099c	YKR028c (11)	unknown	45.8 (1-844)/1058	60.0 (364-1442)/1814
YIL084c	YKR021c (11)	unknown	37.6 (4-932)/1006	46.4 (7-1946)/3138
YIL083c	YKR019c (11)	unknown	26.7 (38-604)/604	64.6 (1265-3800)/1332
YIL082c	YKR018c (11)	unknown	66.0 (1-730)/731	53.7 (233-1986)/1986
YIL079c	YKR013c (11)	unknown	47.5 (1-299)/299	61.4 (415-789)/897
YIL078c	YKR013c (11)	unknown	67.3 (15-651)/681	38.0 (1295-1711)/2843
YIL076c	YKR010c (11)	unknown	16.1 (1-772)/1089	33.7 (2103-3317)/3387
YIL045c	SDH1 (11)	succinate dehydrogenase flavoprotein precursor	83.5 (1-634)/634	78.6 (670-1766)/1902
YIL034c (KAR2)	SSA1 (1)	nuclear fusion protein KAR2 precursor	63.5 (50-663)/682	67.0 (156-1962)/2046
YIL034c (SSC1)	YEL030c (5)	heat shock protein	82.6 (17-642)/654	75.8 (205-1809)/1962
YIR047c (ANB1)	YEL054c (5)	translation initiation factor	90.4 (2-157)/157	91.4 (1-465)/471
YIR048c (CYC1)	YEL039c (5)	cytochrome isoform 1	85.8 (2-107)/109	81.9 (113-323)/327
YIR049c (UTR1)	YEL041c (5)	UTR1 protein	57.0 (304-509)/530	63.8 (419-1392)/1590
YIR051c (OSM1)	YEL047c (5)	involved in osmotic regulation	63.5 (36-499)/501	63.7 (218-1469)/1500
YIR066c (TOR1)	TOR2 (11)	phosphatidylinositol kinase	68.0 (82-2470)/2470	67.2 (2786-7403)/7410
YIR103c (URK8)	URA7 (2)	CTP synthase	79.0 (1-562)/564	71.7 (146-1631)/1692
YIR155c	NOX00 (14)	similar to aryl-alcohol dehydrogenase	89.9 (1-288)/288	87.7 (1-389)/384
YIR156c	NOX95 (14)	similar to thiamine-repressed nat-1	99.8 (1-340)/340	98.4 (568-3311)/1020
YIL274c	YIL276c	similar to α -glucosidase MAL35 (546087)	66.3 (11-587)/599	62.8 (199-1767)/1787
YIL219c	YIL216c	similar to hexose transport protein LGT7	65.2 (33-567)/567	66.3 (226-1685)/1701
YIL078c	YIL078c	unknown	66.7 (152-298)/299	66.2 (551-861)/897
YIL052c (TDH1)	Y8009c (TDH2)	glyceraldehyde-3-phosphate dehydrogenase	65.0 (1-331)/331	92.4 (1-998)/998
YIL038c	YIL037c	unknown	36.3 (5-218)/219	34.0 (295-640)/657

^aWhere known, chromosomal location is indicated in parenthesis.

^bFunction of genes on chromosome X, when available, or else function of their homologues on other chromosomes.

^cNumbers indicate % of aa identity, boundaries of aa comparison (in brackets) and size of the ORF on chromosome X (number after dash).

^dSame as above, but in nt.

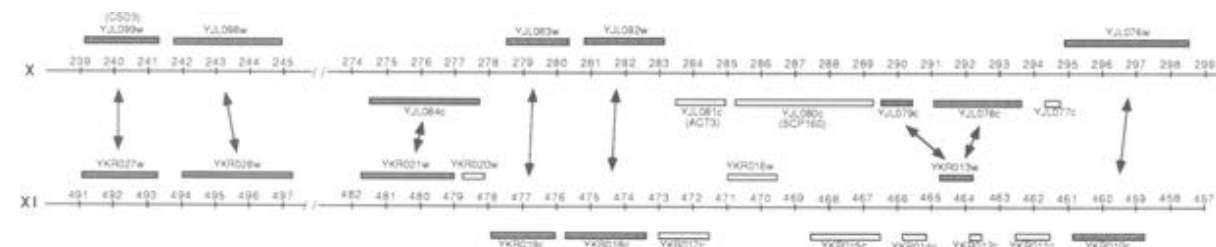


Fig. 4. Physical comparison of the location of genes and systemic segments on chromosome X with that of their counterparts on other chromosomes. The precise position of the genes was deduced from the present sequence and re-drawn to scale coordinates are in kb. Elements above and below the scale belong to the Watson and the Crick strands, respectively. Shaded boxes represent the ORFs with a counterpart on the other chromosome. On the whole, physical distance (and the structures located therein) between any two ORFs on the same systemic segment is not respected on chromosomes other than X. Exceptions are the consecutive ORFs YIL099c (C101) and YIL098c on chromosome X and their homologues YKR027c and YKR028c on chromosome XI, the consecutive ORFs YIL083c and YIL082c on chromosome X and their homologues YKR005c and YKR006c.

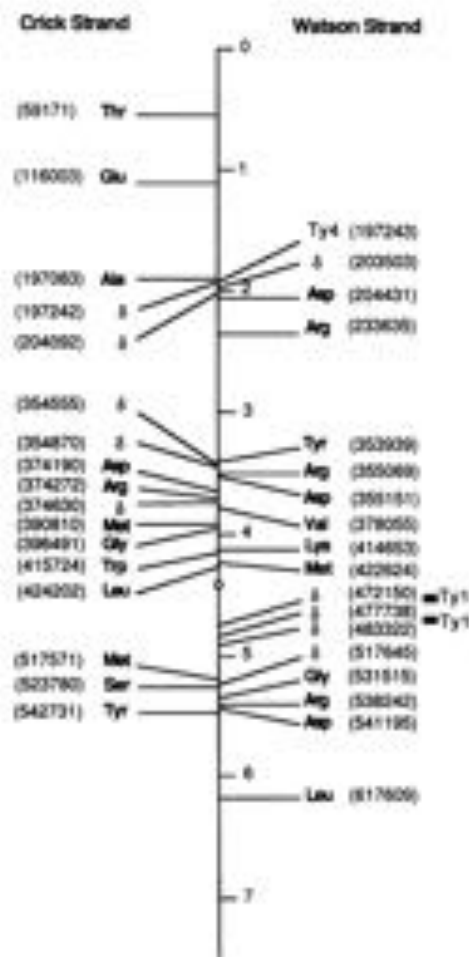


Fig. 5. Position of tRNA genes, Ty sequences and LTRs on chromosome X. The positions were drawn to scale relative to the complete sequence. Elements on the Watson and Crick strands are displayed on the right- and left-hand side, respectively. Only the 3' coordinate is given.



Fig. 6. A clover-leaf structure of yeast tRNA^{Met} on chromosome X (422 624–422 696). All canonical bases are indicated by asterisks. Mismatched base pairs in the stems are boxed. The shadowed nucleotides are the anticodon.

X display an intron 3' to the anticodon sequence, as previously observed. These include two tRNA^{Tyr} with an intron of 14 nt, one of the two tRNA^{Met} with a 19-nt intron and the unique tRNA^{Leu} with an intron of ~29 nt. Its exact size is difficult to assess because base pairing is possible between several short sequences in the anticodon stem, creating an extra arm of variable length.

The entire chromosome X sequence was scanned in parallel for the presence of complete Ty elements or solo remnants or LTR thereof. As shown in Figure 5, several of these have been found. One complete Ty4 is present at position 197243–203468 and two complete Ty1 at position 472150–483659. The two elements are in tandem and share a central δ element. In addition, several solo LTRs are observed. As reported, with the exception of Ty1 these elements are located in the vicinity of tRNA sequences. However, this association seems to be rather loose and, besides, it involves partners located on either strand relative to one another.

Comparison of the physical and genetic maps of the chromosome X

The genetic map of chromosome X includes 60 genes or markers, of which 48 were mapped in a linear array and 12 remained unmapped (Mortimer *et al.*, 1995). Figure 7 shows a comparison of this map with the physical map deduced from the complete nucleotide sequence. Contrary to what has been reported for chromosome XI (Dajon *et al.*, 1994), no gross translocation or inversion was observed here. On the whole, the intergenic distance on the genetic map is roughly proportional to the physical distance, indicative of a relatively uniform recombination frequency over chromosome X. However, closer examination reveals some interesting discrepancies. First, genetic mapping has assigned the previously sequenced *CYR1* gene (alias *CDC35*, *HSR1*, *SRA4* and *TSM0183*), encoding adenylyl cyclase, to a site indistinguishable from that of *met2*. This assignment is clearly incorrect, as the sequence data shows that this gene is in fact located on the left arm of the chromosome, close to the centromere. Second, marked differences are observed in map distances, the ratio between genetic and physical map distances ranging from 0.02 cM per kb for the *TDH2/met1* marker pair, to 0.84 and 4.74 cM per kb for the *met1/lys1* and *lys1/arg1* pairs, respectively. The relatively high frequency of recombination observed in these latter intervals strongly suggests the existence of preferred sites for the initiation of meiotic recombination, similar to those found in the *arg4* region on chromosome VIII (Nicolas *et al.*, 1989; Sun *et al.*, 1989) and the *MAT/ho4* region on chromosome III (Jacquet *et al.*, 1991). It is interesting to note that these intervals of high recombination frequencies in chromosome X appear to coincide with the sharp peak in the G+C content in the right arm of the chromosome (Figure 2).

In all, 31 of the mapped and one, tRNA^{Met}, of the unmapped could be unambiguously assigned to an ORF or a tRNA gene on the basis of sequence comparison. A total of 28 loci cannot at present be attributed to specific ORFs on the physical map of chromosome X.

Discussion

The various elements of the chromosome X sequence referred to above are depicted in Figure 8. The present

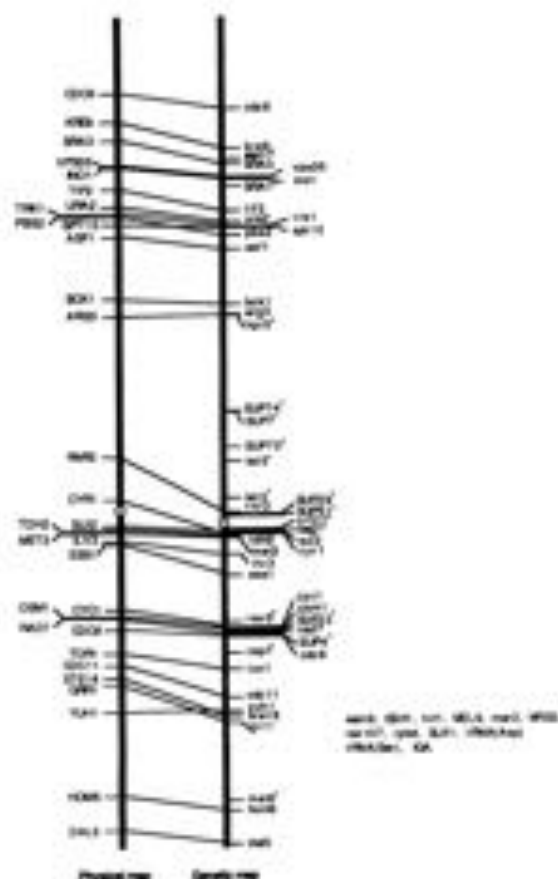


Fig. 7. Comparison of the genetic and physical maps of yeast chromosome X. The genetic map is re-drawn from Mortimer (Mortimer *et al.*, 1995). The unassigned genes or markers are listed on the right. The physical map deduced from this work has been drawn to scale. The circle indicates the position of the centromere. Genes or markers for which no corresponding ORF has been identified on the physical map are indicated by an asterisk.

report brings the number of completely sequenced chromosomes from the yeast *S.cerevisiae* to nine, chromosome X ranking second in this series by virtue of its size. Thus, nearly 40% of the *S.cerevisiae* genome sequence is now accessible to analysis, availability of the whole sequence being anticipated for 1997. The sequence of chromosome X has been established in S288C, a *S.cerevisiae* strain chosen by all members of the European Union sequencing consortium led by André Goffeau. While the study of this sequence reveals no features that are specific for chromosome X, it corroborates several observations made with the previously sequenced chromosomes.

Taking into account only those ORFs whose characteristics, such as size, CAI and disposition leave no doubt as to their existence as real genes, a minimal density of one gene per 2000 nt can be estimated. All these genes are regularly spaced along the chromosome, with no predilection for either strand. Following translation and comparison of the deduced amino acid sequence with database entries, the products of these ORFs can be categorized as follows: (i) 102 proteins previously identified in *S.cerevisiae* and encoded by genes already assigned to chromosome X; (ii) 16 proteins with strong similarity,

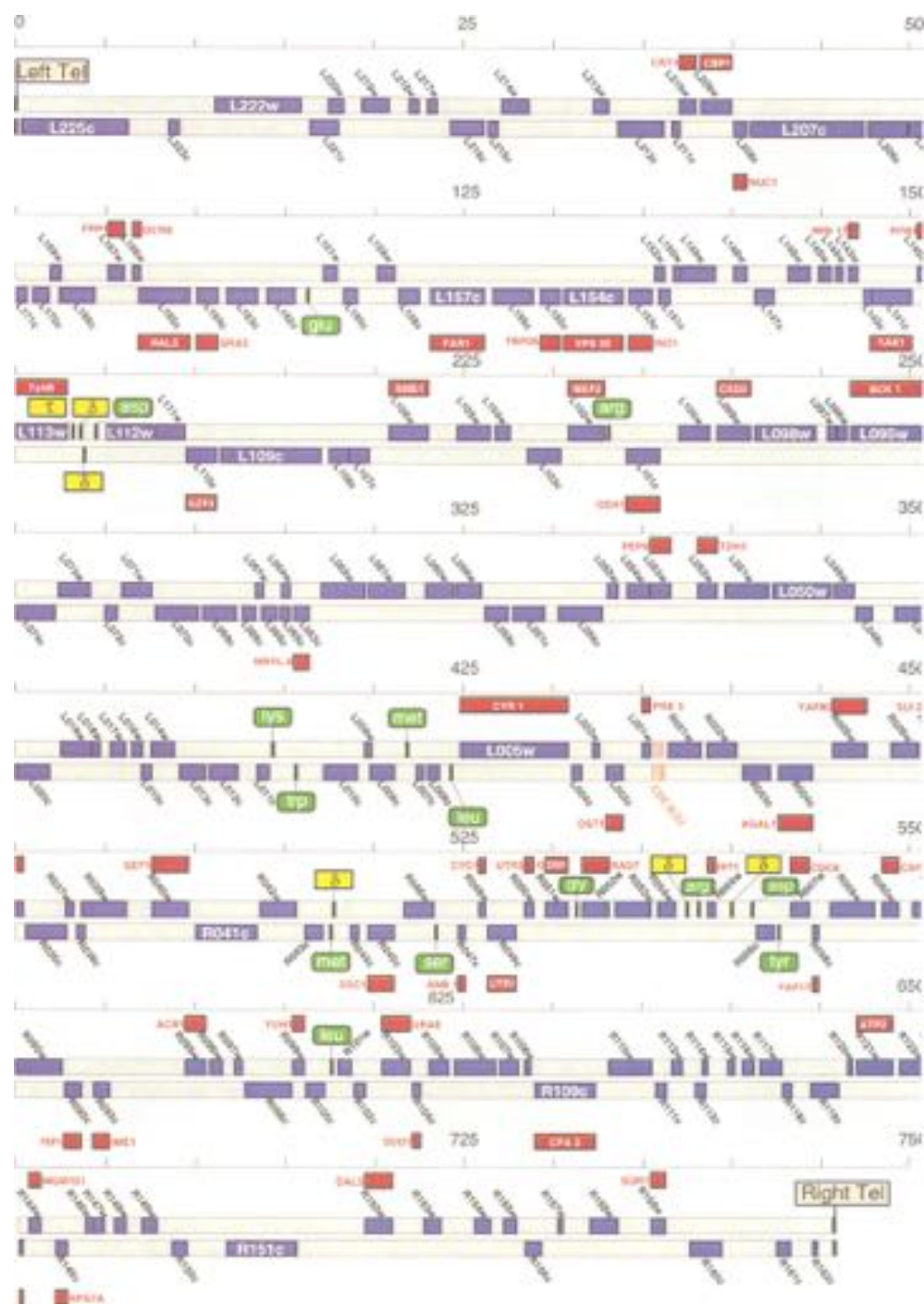
or even near identity, to known *S.cerevisiae* proteins, but whose coding gene has not previously been shown to reside on chromosome X; (iii) 22 proteins with a FastA score much greater than 200—equal to at least half the self-score, i.e. the score obtained when the protein is compared with itself. Such high scores can be considered as warranting a realistic hypothesis regarding the function of ORFs in this category; (iv) 35 proteins with a FastA score >200, though lower than half the self-score. A function can also be envisaged in this case, but with more caution; (v) 92 proteins with no significant FastA score but displaying a particular motif signature; (vi) 112 proteins with no match at all in database entries. This last category remains numerically important, since it includes nearly 30% of the ORFs, a proportion that fully vindicates the systematic sequencing approach of the *S.cerevisiae* genome launched in 1989.

Regarding ORFs in categories (iii) and (iv) above, for which a function can be hypothesized, several of the proteins discovered in chromosome X are worth mentioning. For instance, three new genes encoding different subunits of the cytosolic chaperone complex (CCT5, CCT7 and CCT8) have been discovered on chromosome X in addition to CCT3. This brings the number of fully sequenced CCT genes in *S.cerevisiae* to eight. Together with the versatility of yeast versus mouse genetics, availability of these sequence data will undoubtedly promote fine molecular analysis of this important chaperone system. Another remark concerns the discovery of a Cl⁻ channel gene (Huang *et al.*, 1994c) on chromosome X. In this respect, it is both surprising and remarkable that systematic sequencing was required to detect the first Cl⁻ channel ever described in a species as thoroughly studied as *S.cerevisiae*. Here again, availability of the gene and of disruption mutants thereof will permit identification by complementation homologous genes in other species of interest, in particular in plants.

Chromosome X stands out because of the number of tRNA genes (24) it accommodates, capable of transferring 13 different amino acids. However, what is even more remarkable and has so far escaped notice is that folding of these tRNAs according to the clover-leaf model reveals quite a few mismatches in the several stems. This is suggestive of an editing process aiming at correcting some of these mismatches, as reported for various tRNAs from plants (Maréchal-Drouot, 1993). Of course, validation or dismissal of this hypothesis must await analysis at the RNA level.

Duplicated genes are found in chromosome X, as in other *S.cerevisiae* chromosomes. These include both intra- and interchromosomal duplications. Furthermore, actual systemic regions can be recognized in the latter case. The implications are 2-fold, pertaining (i) to the study of the evolution of the yeast genome and (ii) to function analysis, as it is known that disruption of a single gene frequently does not result in any phenotypic alteration. By the same token, a clue to the function of a gene might in some instances be provided by disruption of all the genes belonging to a given family.

To conclude, it must be stressed that this brief account of the sequence analysis of chromosome X cannot cover all the information embedded in the nucleotide sequence



and that many biological analyses will be needed to exploit this mine of information in the years to come.

Materials and methods

Chromosome X DNA

Total yeast DNA was obtained from FY1679, a diploid strain issued from the cross between strains FY23 (MAT α , ara1-32, rpl36J, leu2 Δ , GAL2) and FY77 (MAT α , ara1-32, his AS200, GAL2). FY23 and FY77 are derived from strain S288C and are isogenic with it except for the

markers indicated (Wiseman *et al.*, 1995). The construction of an ordered cosmid library and of an *Eco*RI restriction map have been previously published (Huang *et al.*, 1994a). Overlapping cosmids covering the chromosome X coding were distributed within a consortium of 15 laboratories. The telomeres and subtelomeric regions were cloned in vector pIL51, as described by Louis and Borts (1995).

Determination, assembly and analysis of the sequence

Sequencing strategies and methods varied among the 15 collaborating laboratories (Table V). Sequence assembly in the single contracting laboratories was performed by a variety of software program packages.

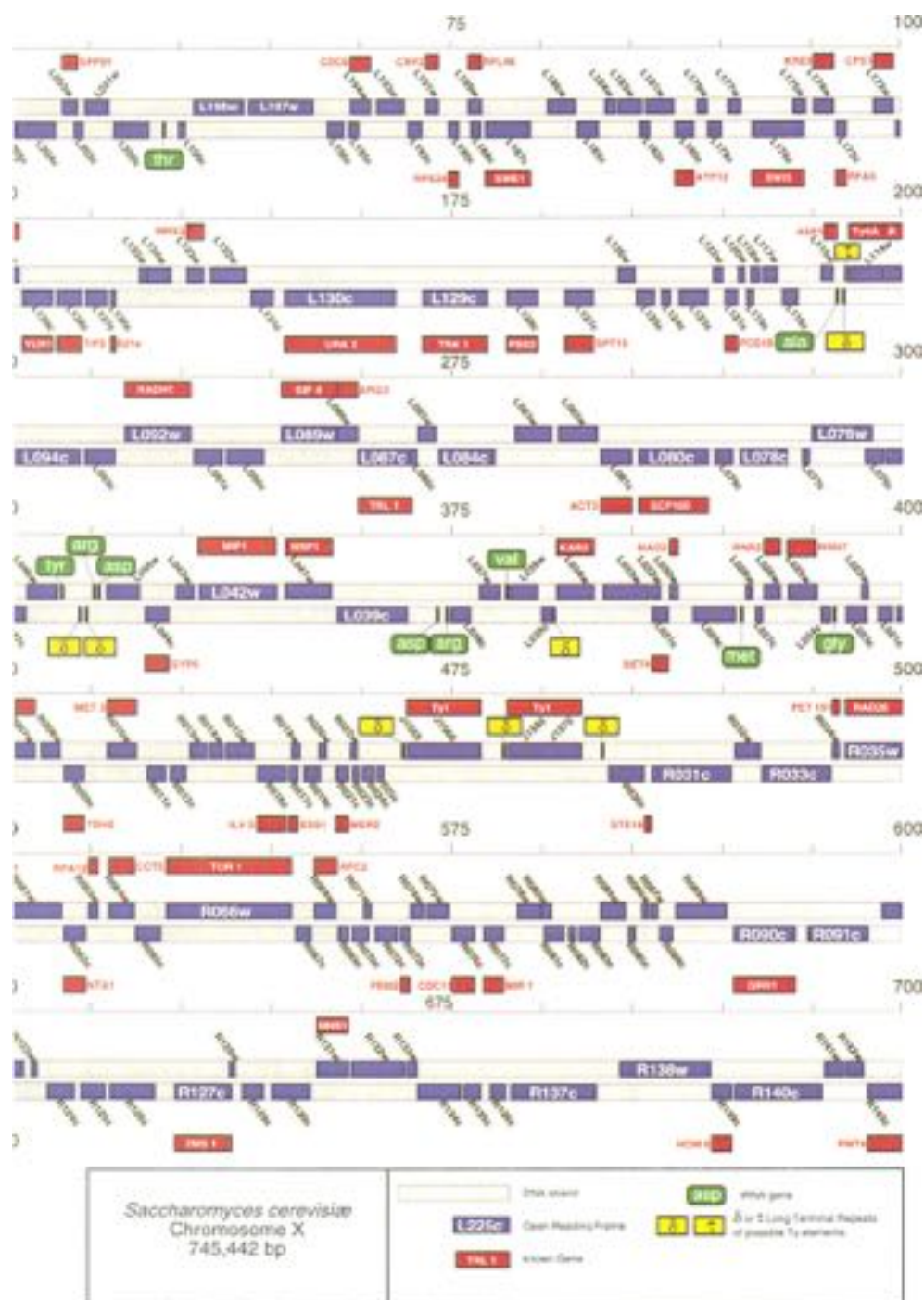


Fig. 8. Chromosome X map deduced from the complete sequence. The chromosome and its constitutive elements are drawn to scale. The top bar represents the Watson strand oriented 5' to 3' from left to right, the bottom bar the Crick strand. The conserved elements of the centromere are designated as CDE I, II and III. ORFs on the left and right arm are designated by the letters L and R, respectively, before their number (numbering is in increasing order from the centromere). Full designations, in accordance with the official ORF nomenclature, are obtained by adding again the letters Y (for yeast) and X (for chromosome X) at the beginning, and w (Watson) or c (Crick) at the end.

The telomeres were cloned in Oxford. The left telomere was sequenced in one of 15 laboratories. The right telomere and the PCR fragment filling the gap were sequenced in Berlin. Completed contigs submitted to MIPS were stored in a data library and assembled using the GCG software package 7.2 for the NAX (Deveraux *et al.*, 1984). The names

and position of genetic elements have been deduced from the sequence using the following principles: (i) all possible in-frame open reading frame pairs were detected using specially defined patterns (Fondra *et al.*, 1994); (ii) ORFs occurring in all possible frames were listed. ORFs containing at least 99 contiguous sense codons following an ATG and

Table V. Methods used by each of the collaborating laboratories

Whole cosmid Shotgun	Restricted Fragments		
	Shotgun	TN/OD	Nested deletions
Louvain (M)	Gemboux (M)	Darmstadt (M)	München (A)
Heidelberg (M)	Amsterdam (A)	Frankfurt (A)	Copenhagen (A)
Konstanz (M)			Düsseldorf (A)
Paris (A)			Ghent (A)
Gif (A)			Heidelberg (M)
Bonn (A)			

M, manual methods; A, automated methods.

those containing 20–98 codons were retained for further analysis, in both cases provided they were not entirely contained within a longer ORF on either DNA strand. Searches for similarity of the deduced protein sequences to entries in the databanks were performed by FastA (Pearson and Lipman, 1988) in the Protein Sequence Database of PIR International (release 44) and other databanks. Protein signatures were detected using the PROSITE dictionary (release 11.1) (Battey, 1989). ORFs were assigned probable functions when the alignments from FastA searches showed significant similarity and/or protein signatures were apparent, whereas FastA scores <200 were considered insufficient to confidently assign function. The complete sequence was also searched for tRNA genes ('tscan') (Fichant and Burks, 1991), consensus and telomere consensus elements and for 5', 6' or 7' elements by comparison with a data set of such elements previously characterized in yeast. Compositional analyses of the chromosome were performed using the X11 program package (C. Marck, unpublished results). For calculations of CAI and GC content of ORFs, the algorithm CODONS (Lloyd and Sharp, 1992) was used.

Sequence verifications and quality controls

All sequences submitted by collaborating laboratories to the Maxamill Institute for Protein Sequences (MIPS) data library were subjected to quality controls. The procedure was comprised of three major steps. First, the strategy of each contractor was checked by the coordinator to pinpoint possible weak points and request the sequencers to review their electropherograms to assess the quality of their reads in these less documented regions. Second, once cosmid sequences had been entered in the database, the match between the overlaps was held to provide an assessment of the respective quality of the neighbouring partial sequences. Third, each of the cosmids that had been distributed to the contractors for sequencing was shotgunned, size-selected to ~300–500 bp and cloned in plasmid vector, the size of the inserts ensuring that sequencing with the universal forward and reverse primers would provide a 300–600 double-stranded sequence. The subclones from each cosmid were sent with coded names to a different sequencer. The double-stranded part of each sequence was then sent to MIPS and compared with the initial sequence. The number of verification sequences per cosmid clone (averaging 15–30) varied according to the quality of the initial sequencing as deduced from alignment within the overlaps. Any discrepancy detected between overlapping partial sequences or between the sequence initially submitted and the verification sequence was addressed as follows. A stretch of 20 bp including the discrepancy, but not centering on it, was pointed out to each party for reviewing and re-submission to MIPS, whether modified or not. This procedure was sufficient to remove most discrepancies, as one party usually provided a revised sequence matching the other's. Resistant cases were dealt with by requesting both parties to send the electropherograms corresponding to the conflicting sequences to the coordinator, who made a decision and requested resequencing if necessary.

The sequence data reported are available through <http://mips.biochem.mpg.de/yeast>

Acknowledgements

We wish to thank B.Dujon for fruitful discussions and for help with the gene density and G+C composition plots, and G.Le Prevost for secretarial assistance. The laboratory consortium operating under contracts with the European Commission was initiated and organized by A.Giffoux. This study is part of the second phase of the European Yeast Genome Sequencing Project carried out under the administrative coordination of

A.Vissarotti (DG-XII) and the Université Catholique de Louvain, and under the scientific responsibility of F.Galbert as DNA coordinator. This work was supported by the European Commission under the BRIDGE and Biotech programmes, the Groupe de Recherche et d'Etudes sur le Génome (GREG) and the Centre National de la Recherche Scientifique (CNRS) (PR), the Wellcome Trust (UK), the Région Wallonne and the Fonds National de la Recherche Scientifique (BE), the Bundesminister für Forschung und Technologie (DE) and the Ministry of Industry and Technology (GR).

References

- Battey, J. (1989) EMBL Bioinformatics Technical Document 4. EMBL, Heidelberg, Germany.
- Barell, R.G. et al. (1994) Sequence of *S. cerevisiae* chromosome IX. <http://www.sanger.ac.uk/~yeastpub/seq/w/sequencing.html>
- Bussay, H. et al. (1993) The nucleotide sequence of chromosome I from *Saccharomyces cerevisiae*. *Proc. Natl. Acad. Sci. USA*, **92**, 3809–3813.
- Deliboudier, A., Gopali, V., Bocan, A.M., Couvrot, F., Perra, J., Rantopuu, J. and Jacq, C. (1989) Site-specific DNA endonuclease and RNA maturation activities of two homologous intron-encoded proteins from yeast mitochondria. *Cell*, **56**, 431–441.
- Deveraux, J., Haeberli, P. and Smithies, D. (1984) A comprehensive set of sequence analysis programs for the VAX. *Nucleic Acids Res.*, **12**, 387–395.
- Dietrich, F.S. et al. (1994) Sequence of *S. cerevisiae* chromosome V. <http://openly.mips.biochem.mpg.de/mips/yeast/ch5>
- Dujon, B. et al. (1994) Complete DNA sequence of yeast chromosome XI. *Nature*, **368**, 371–378.
- Feldmann, H. et al. (1994) Complete DNA sequence of yeast chromosome II. *EMBO J.*, **13**, 5793–5809.
- Fichant, G.A. and Burks, C. (1991) Identifying potential tRNA genes in genomic DNA sequences. *J. Mol. Biol.*, **220**, 659–671.
- Fondral, C. and Katsouroulis, A. (1994) Approaching the function of new genes by detection of their potential upstream activation sequences in *Saccharomyces cerevisiae*: application to chromosome III. *Curr. Genet.*, **25**, 396–406.
- Galbert, F., Mandart, E., Fissore, F., Toffais, P. and Charney, P. (1979) Nucleotide sequence of the hepatitis B virus genome (subtype ayw) cloned in *E. coli*. *Nature*, **281**, 646–650.
- Hemrika, W., Barden, J.A. and Gissel, L.A. (1993) A region of the C-terminal part of the 11-kDa subunit of ubiquitin-cytochrome-c oxidoreductase of the yeast *Saccharomyces cerevisiae* contributes to the structure of the Q-sorts reaction domain. *Eur. J. Biochem.*, **215**, 601–609.
- Haier, P., Pridmore, D., Hegeman, J.H., Thomas, M., Davis, R.W. and Philippon, P. (1992) Functional selection and analysis of yeast centromeric DNA. *Cell*, **62**, 913–921.
- Honebusch, A.G. and Liebman, S.W. (1991) Protein synthesis and translational control in *Saccharomyces cerevisiae*. In Broach, J.R. et al. (eds), *The Molecular Biology of the Yeast Saccharomyces*. Cold Spring Harbor Laboratory Press, Plainview, NY, pp. 627–715.
- Huang, M.F., Chua, J.C., Thierry, A., Dujon, B. and Galbert, F. (1994a) Construction of a cosmid contig and of an EcoRI restriction map of yeast chromosome X. *DNA Sequence*, **4**, 263–300.
- Huang, M.F., Manos, V., Chua, J.C. and Galbert, F. (1994b) Revised nucleotide sequence of the COR region of yeast *S. cerevisiae* chromosome X. *Trust*, **10**, 811–818.

- Huang, M.E., Chan, J.C. and Gilbert, F. (1994a) A voltage-gated chloride channel in the yeast *Saccharomyces cerevisiae*. *J. Mol. Biol.*, **242**, 595-598.
- Huang, M.E., Chan, J.C. and Gilbert, F. (1995) Analysis of a 42.5 kb DNA sequence of chromosome X reveals three tRNA genes and 14 new open reading frames including a gene most probably belonging to the family of ubiquitin-protein ligase. *Nucl. Acids Res.*, **23**, 775-781.
- Jacquet, M., Babler, J.M., Iborra, F., Francini-Guiffard, M.C. and Soustelle, C. (1991) The MAT locus revisited within a 9.8 kb fragment of chromosome III containing BUD3 and two new open reading frames. *Yeast*, **7**, 881-888.
- Johnson, M. et al. (1994) Complete nucleotide sequence of *Saccharomyces cerevisiae* chromosome VIII. *Science*, **265**, 2077-2082.
- Kerchov, K.A., Jones, W.J., Sellers, A.J., Kaczmarek, L.K. and Goldstein, S.A.N. (1995) A new family of outwardly rectifying potassium channel proteins with two pore domains in tandem. *Nature*, **376**, 690-695.
- Korak, M. (1987) An analysis of 5'-noncoding sequences from 699 vertebrate messenger RNAs. *Nucleic Acids Res.*, **15**, 8125-8148.
- Liby, L.A.T. and Sharp, P.M. (1992) CODONS: A microcomputer program for codon usage analysis. *J. Mol. Evol.*, **33**, 239-240.
- Louis, E.J. and Born, R.H. (1995) A complete set of marked telomers in *Saccharomyces cerevisiae* for physical mapping and cloning. *Genetics*, **139**, 125-136.
- Louis, E.J. and Haber, J.E. (1991) Evolutionarily recent transfer of a group I mitochondrial intron to telomere regions in *Saccharomyces cerevisiae*. *Curr. Genet.*, **20**, 411-415.
- Louis, E.J., Naumova, E.S., Lee, A., Naumov, G. and Haber, J.E. (1994) The chromosome end in yeast: its mosaic nature and influence on recombinational dynamics. *Genetics*, **136**, 789-802.
- Marichal-Drouard, L., Ramarompong, D., Couet, A., Weill, H. and Dietrich, A. (1993) Editing correct mispairing in the acceptor stem of bean and potato mitochondrial phenylalanine transfer RNAs. *Nucleic Acids Res.*, **21**, 4909-4914.
- Mingq, T., Witzel, A. and Zimmermann, F.K. (1994a) Sequence and function analysis of a 9.46 kb fragment of *Saccharomyces cerevisiae* chromosome X. *Yeast*, **10**, 965-973.
- Mingq, T., Böles, E., Schaaff-Gerstenschläger, J., Schmitt, S. and Zimmermann, F.K. (1994b) Sequence and function analysis of a 9.74 kb fragment of *Saccharomyces cerevisiae* chromosome X including the BCK1 gene. *Yeast*, **10**, 1401-1408.
- Mingq, T., Schaaff-Gerstenschläger, J., Chabwaits, N., Baur, A., Böles, E., Fournier, C., Schmitt, S., Vollen, C., Wilhelm, N. and Zimmermann, F.K. (1995) Sequence analysis of a 33.1 kb fragment from the left arm of *Saccharomyces cerevisiae* chromosome X, including putative proteins with leucine zippers, a fungal Zn(II)-Cys6 binuclear cluster domain and a putative alpha2-SCB-binding site. *Yeast*, **11**, 681-689.
- Mortimer, R.K., Cherry, J.M., Dietrich, F.S., Riles, L., Olson, M.S. and Bonstein, D. (1995) Genetic map of *Saccharomyces cerevisiae*. Edition 17. <http://genome-www.stanford.edu/saccharbioedition17.html>
- Murakami, Y. et al. (1995) Analysis of the nucleotide sequence of chromosome VI from *Saccharomyces cerevisiae*. *Nature Genet.*, **10**, 261-268.
- Navarre, C., Ghidoui, M., Leterme, S., Ferroud, C., Dufour, J.P. and Goffeau, A. (1992) Purification and complete sequence of a small proteinolipid associated with the plasma membrane H⁺-ATPase of *Saccharomyces cerevisiae*. *J. Biol. Chem.*, **267**, 6425-6428.
- Navarre, C., Catty, P., Leterme, S., Dietrich, F. and Goffeau, A. (1994) Two distinct genes encode small isoprenolipids affecting plasma membrane H⁺-ATPase activity of *Saccharomyces cerevisiae*. *J. Biol. Chem.*, **269**, 21262-21268.
- Nicolaus, A., Tracy, D., Schultes, N.P. and Szostak, J.W. (1989) An initiation site for meiotic gene conversion in the yeast *Saccharomyces cerevisiae*. *Nature*, **333**, 87-90.
- Olivier, S. et al. (1992) The complete DNA sequence of yeast chromosome III. *Nature*, **357**, 38-46.
- Pearson, W.R. and Lipman, D.J. (1988) Improved tools for biological sequence comparison. *Proc. Natl. Acad. Sci. USA*, **85**, 2444-2448.
- Pryde, F.E., Hackle, T.C. and Louis, E.J. (1995) Sequence analysis of the right end of chromosome XV in *Saccharomyces cerevisiae*: An insight into the structural and functional significance of sub-telomeric repeat sequences. *Yeast*, **11**, 371-382.
- Porelle, B., Couet, F. and Goffeau, A. (1994) The sequence of a 36 kb segment on the left arm of yeast chromosome X identifies 24 open reading frames including NUC1, PRP21(PPP1), CDC6, CRP2, the gene for S24, a homologue to the acetonase gene ACO1 and two homologues to chromosome III genes. *Yeast*, **10**, 1235-1249.
- Pitt, J., Puglisi, J.D., Florentz, C. and Gagli, R. (1993) Additive, cooperative and anti-cooperative effects between identity nucleotides of a tRNA. *EMBO J.*, **12**, 2949-2953.
- Rasmussen, S.W. (1995) A region of yeast chromosome X includes the SME1, MEF2, GSH1 and CSD3 genes, a TCP-1 related gene, an open reading frame similar to the DAL80 gene, and a tRNA^{Pro}. *Yeast*, **11**, 873-883.
- Sharp, P.M. and Li, W.H. (1987) The codon adaptation index — a measure of directional synonymous codon usage bias, and its potential applications. *Nucleic Acids Res.*, **15**, 1281-1295.
- Sun, H., Tracy, D., Schultes, N.P. and Szostak, J.W. (1989) Double stranded breaks at an initiation site for meiotic gene conversions. *Nature*, **333**, 87-90.
- Vanderbol, M., Durand, P., Balle, P.-A., Dion, C., Portetello, D. and Hilger, F. (1994) Sequence analysis of a 40.2 kb DNA fragment located near the left telomere of yeast chromosome X. *Yeast*, **10**, 1657-1662.
- Vanderbol, M., Durand, P., Dion, C., Portetello, D. and Hilger, F. (1995) Sequence of a 17.1 kb DNA fragment from chromosome X of *Saccharomyces cerevisiae* includes the mitochondrial ribosomal protein L8. *Yeast*, **11**, 57-60.
- Winston, F., Dollard, C. and Ricupero-Hovavest, S.L. (1995) Construction of a set of convenient *Saccharomyces cerevisiae* strains that are isogenic to S288C. *Yeast*, **11**, 53-55.
- Zaputski, M., Babinska, B., Gromadka, R., Mijalowski, A., Sulicka, J. and Herbert, C.J. (1995) The sequence of 24.3 kb from chromosome X reveals 5 complete open reading frames all of which correspond to new genes, and a tandem insertion of a Ty1 transposon. *Yeast*, **11**, 1179-1186.

Received on November 3, 1995; revised on January 5, 1996