ELSEVIER

# Using a Euclid distance discriminant method to find protein coding genes in the yeast genome

Chun-Ting Zhang [a,*], Ju Wang [a], Ren Zhang [b]

[a] *Department of Physics*, *Tianjin University*, *Tianjin 300072*, *China*
[b] *Department of Epidemiology and Biostatistics*, *Tianjin Cancer Institute and Hospital*, *Tianjin 300060*, *China*

## Abstract

The Euclid distance discriminant method is used to find protein coding genes in the yeast genome, based on the single nucleotide frequencies at three codon positions in the ORFs. The method is extremely simple and may be extended to find genes in prokaryotic genomes or eukaryotic genomes with less introns. Six-fold cross-validation tests have demonstrated that the accuracy of the algorithm is better than 93%. Based on this, it is found that the total number of protein coding genes in the yeast genome is less than or equal to 5579 only, about 3.8–7.0% less than 5800–6000, which is currently widely accepted. The base compositions at three codon positions are analyzed in details using a graphic method. The result shows that the preference codons adopted by yeast genes are of the $R\bar{G}W$ type, where R, $\bar{G}$ and W indicate the bases of purine, non-G and A/T, whereas the 'codons' in the intergenic sequences are of the form NNN, where N denotes any base. This fact constitutes the basis of the algorithm to distinguish between coding and non-coding ORFs in the yeast genome. The names of putative non-coding ORFs are listed here in detail. © 2002 Elsevier Science Ltd. All rights reserved.

*Keywords:* Gene-finding; Yeast genome; Z curve; Euclid distance discriminant method

## 1. Introduction

The budding yeast *Saccharomyces cerevisiae* is an important model organism for the Human Genome Project. Due to the efforts of more than 600 scientists worldwide, the first sequenced genome of an eukaryotic organism, *S. cerevisiae*, has been completed (Oliver et al., 1992; Dujon et al., 1994, 1997; Bowman et al., 1997; Feldmann et al., 1994; Galibert et al., 1996; Johnston et al., 1994, 1997; Bussey et al., 1995, 1997; Murakami et al., 1995; Churcher et al., 1997; Dietrich et al., 1997; Jacq et al., 1997; Philippsen et al., 1997; Tettelin et al., 1997). Although this is a great scientific achievement,

much work remains to be done. The completion of the Yeast Genome Project may be deemed as the first step in a 'Long March' towards understanding the genetic secret of this relatively simple organism. It is necessary to clarify functions of genes and relationships of them. However, to clarify the number of genes is even a more critical task at present. The number of protein coding genes in the yeast genome was estimated to be 5800–6000 (Goffeau et al., 1996; Winzeler and Davis, 1997; Mewes et al., 1997), which is currently widely accepted. On the contrary, another group estimated recently that the number should be less than 4700 (Mackiewicz et al., 1999). The results are obviously controversial.

Historically, the codingness of an ORF or a fragment of DNA sequence in the yeast genome was described by using the Codon Bias Index (CBI) (Benetzen and Benjamin, 1982) or the Codon Adaptation Index (CAI)

---

* Corresponding author. Tel.: + 86-22-2740-2987; fax: + 86-22-2335-8329.

*E-mail address:* ctzhang@tju.edu.cn (C.-T. Zhang).

(Sharp and Li, 1987). Although these indices were used widely (Dujon et al., 1994), the coding properties of a coding sequence are not sufficiently reflected by them. For example, some ORFs shorter than 150 codons with $CAI < 0.11$ have identified phenotypes (Mackiewicz et al., 1999). During the past decade, numerous advanced gene-finding algorithms have been developed. See, e.g., a recent review paper written by Fickett (1996). A set of 745 sequences in the yeast genome was selected to evaluate the gene-finding algorithm based on a correspondence analysis (Quentin et al., 1999). Genes in the same set were also predicted using the GeneMark program (Borodovsky and McIninch, 1993) (see the discussion of Quentin et al. (1999)). In this paper an extremely simple gene-finding algorithm based on the Euclid distance discriminant method is proposed. The algorithm utilizes a graphic approach to explore the difference between coding and non-coding sequences. In addition, this simple gene-finding algorithm is useful because it may be complementary with other existing methods. Therefore, by joining them, more reliable gene recognition results could be expected. Based on the algorithm, a new index is proposed to describe the codingness of yeast ORFs, which may be an appropriate complement to CBI or CAI, which are already widely used.

## 2. Materials and methods

### 2.1. The database

The *S. cerevisiae* genome DNA sequences were obtained from a CD-ROM distributed from the Munich Information Centre for Protein Sequences (MIPS), released in 1997. The data for classification of ORFs in the yeast genome were downloaded from http://speedy.mips.biochem.mpg.de, release, September 27, 1999 (Mewes et al., 1999) (the database is referred to as MIPS database, hereafter). In the MIPS database, all the ORFs are classified into six classes, which correspond to known proteins, strong similarity to known proteins, similarity or weak similarity to known proteins, similarity to unknown proteins, no similarity and questionable ORFs, respectively. The 1st, 2nd, 3rd, 4th, 5th and 6th classes include 3199 (18), 248, 869, 789 (1), 805 and 447 (8) entries, respectively, where the figures in the parentheses indicate the numbers of ORFs in the mitochondrial genome. The mitochondrial ORFs are excluded here since the samples are too few to have statistical significance. So in each of the six classes, 3181, 248, 869, 788, 805 and 439 ORFs are contained, respectively.

### 2.2. The gene-finding algorithm

The gene-finding algorithm presented in this paper is based on the differences of single nucleotide frequencies at the three codon positions between protein coding ORFs and non-coding ones. Suppose that the occurrence frequencies of the bases A, C, G and T at the 1st, 2nd and 3rd codon positions in an ORF are denoted by $a_i$, $c_i$, $g_i$ and $t_i$, respectively, where $i = 1, 2, 3$. Since $a_i + c_i + g_i + t_i = 1$, the four real numbers $a_i$, $c_i$, $g_i$ and $t_i$ may be mapped onto a point $P_i$ in a three-dimensional space $V_i$. The coordinates $x_i$, $y_i$ and $z_i$ of $P_i$ are determined by the so-called *Z-transform* for DNA sequence, which transforms the nucleotide frequencies into a three-dimensional curve, the Z curve (Zhang and Zhang, 1991)

$$\begin{cases} x_i = (a_i + g_i) - (c_i + t_i), \\ y_i = (a_i + c_i) - (g_i + t_i), \quad i = 1, 2, 3. \\ z_i = (a_i + t_i) - (g_i + c_i), \end{cases} \quad (1)$$

Define $V = V_1 \oplus V_2 \oplus V_3$, i.e. the nine-dimensional space V is the direct sum of the subspaces $V_1$, $V_2$ and $V_3$. Denoting the nine bases of V by $u_1 - u_9$, we define

$$\begin{cases} u_1 = x_1, & u_2 = y_1, & u_3 = z_1, \\ u_4 = x_2, & u_5 = y_2, & u_6 = z_2, \\ u_7 = x_3, & u_8 = y_3, & u_9 = z_3. \end{cases} \quad (2)$$

Therefore, each coding ORF or non-coding DNA sequence is represented by a point or a vector, respectively, in the nine-dimensional space V.

To complete the algorithm in a computer, usually a training set of samples (ORFs) is needed. The training set consists of two parts: one includes the positive samples composed of true protein coding genes, whereas the other includes negative samples composed of non-coding DNA sequences. Suppose that there are $N$ samples in each part. In the positive samples the $i$-th true coding ORF is described by a vector $(u_{i,1}^1, u_{i,2}^1, \ldots, u_{i,9}^1)^{\mathrm{T}}$, where $u_{i,s}^1$ are the $s$-component of the vector ($s = 1, 2, \ldots, 9$), and 'T' indicates a transpose operator for a matrix. Similarly, in the negative samples the $i$-th non-coding DNA sequence is described by a vector $(u_{i,1}^2, u_{i,2}^2, \ldots, u_{i,9}^2)^{\mathrm{T}}$, where $u_{i,s}^2$ are the $s$-component of the vector ($s = 1, 2, \ldots, 9$). The geometric centers for the positive and negative samples in the 9-dimensional space V are denoted by $\bar{U}^1$ and $\bar{U}^2$, respectively, where

$$\bar{U}^1 = (\bar{u}_1^1, \bar{u}_2^1, \ldots, \bar{u}_9^1)^{\mathrm{T}}, \quad \bar{U}^2 = (\bar{u}_1^2, \bar{u}_2^2, \ldots, \bar{u}_9^2)^{\mathrm{T}}, \quad (3)$$

and

$$\bar{u}_s^1 = \frac{1}{N} \sum_{i=1}^{N} u_{i,s}^1, \quad \bar{u}_s^2 = \frac{1}{N} \sum_{i=1}^{N} u_{i,s}^2, \quad s = 1, 2, \ldots, 9. \quad (4)$$

Suppose that a query ORF is indicated by a nine-dimensional vector $\mathbf{U} = (u_1, u_2, \ldots, u_9)^{\mathrm{T}}$. To judge whether this ORF is a true protein coding gene or not, calculate the Euclid distance $d(\mathbf{U}, \bar{\mathbf{U}}^1)$ between $\mathbf{U}$ and $\bar{\mathbf{U}}^1$, and the Euclid distance $d(\mathbf{U}, \bar{\mathbf{U}}^2)$ between $\mathbf{U}$ and $\bar{\mathbf{U}}^2$, where

$$d(\mathbf{U}, \bar{\mathbf{U}}^1) = \left[ \sum_{s=1}^{9} (u_s - \bar{u}_s^1)^2 \right]^{1/2},$$

$$d(\mathbf{U}, \bar{\mathbf{U}}^2) = \left[ \sum_{s=1}^{9} (u_s - \bar{u}_s^2)^2 \right]^{1/2}. \tag{5}$$

A codingness index $\Delta$ is defined as

$$\Delta = d(\mathbf{U}, \bar{\mathbf{U}}^2) - d(\mathbf{U}, \bar{\mathbf{U}}^1) + c, \tag{6}$$

where $c$ is a constant determined by making the false positive rate and false negative rate identical in the training set. If $\Delta > 0$, the query ORF is recognized as coding gene, otherwise, if $\Delta < 0$, the ORF or DNA sequence is recognized as a non-coding one.

## 3. Results and discussions

### 3.1. Definitions of sensitivity, specificity and accuracy

To evaluate the performance of the algorithm, we have to discuss the definitions of the accuracy, sensitivity and specificity. The notations used here are the same as in Burset and Guigo (1996). Denoted by TP the number of coding ORFs that have been correctly predicted as coding, and denoted by FN the number of coding ORFs that have been predicted as non-coding, we define the sensitivity $s_n$ as

$$s_n = \frac{\mathrm{TP}}{\mathrm{TP} + \mathrm{FN}}. \tag{7}$$

That is, $s_n$ is the proportion of coding ORFs that have been correctly predicted as coding. Similarly, denoted by TN the number of intergenic sequences that have been correctly predicted as non-coding, and denoted by FP the number of intergenic sequences that have been predicted as coding, we define the specificity $s_p$ as

$$s_p = \frac{\mathrm{TN}}{\mathrm{TN} + \mathrm{FP}}. \tag{8}$$

Table 1
The accuracy of the algorithm for three different test sets

| Test set | 1 | 2 | 3 |
|---|---|---|---|
| Sensitivity (%) | 93.00 | 93.60 | 94.60 |
| Specificity (%) | 93.10 | 93.40 | 93.00 |
| Accuracy[a] (%) | 93.05 | 93.50 | 93.80 |

[a] Accuracy is defined as the average of the sensitivity and specificity.

That is, $s_p$ is the proportion of intergenic sequences that have been correctly predicted as non-coding. The accuracy is defined as the average of $s_n$ and $s_p$.

### 3.2. Self-consistency and cross-validation tests

How to evaluate a gene-finding algorithm is an important issue. Usually, the accuracy of a gene-finding algorithm is evaluated by the resubstitution and cross-validation tests, respectively. The former reflects the self-consistency, and the latter reflects the extrapolating effectiveness of the algorithm. To evaluate the algorithm, a training set and a test set are needed, which should be independent of one another. In the MIPS database, the first class includes 3181 known genes residing in the 16 yeast chromosomes. Among them, 223 are intron-containing genes and the remaining 2958 are intronless. Randomly divide the 2958 intronless genes into two unequal parts, in which the larger part includes 1958 genes, whereas the smaller includes 1000 genes. The former is served as a training set, whereas the latter is served as a test set. Both the training and test sets should be accompanied by the counterparts of negative samples. We have randomly selected about 6000 intergenic sequences with lengths longer than 300 bp from the 16 yeast chromosomes, and each of them starts with ATG and ends with one of the stop codons. Note that such sequences are unlikely to be ORFs, because there are usually several stop codons within the sequences. We randomly selected 1958 and 1000 intergenic sequences from the above 6000 sequences, which form the training and test sets of the negative samples, respectively. In summary, the training set includes 1958 positive samples (true genes) and 1958 negative samples (intergenic sequences). The test set includes 1000 positive samples (true genes) and 1000 negative samples (intergenic sequences). Using the sequences in the training set, the average vectors $\bar{\mathbf{U}}^1$, $\bar{\mathbf{U}}^2$ and the parameter $c$ (see Eq. (6)) are determined. Using these quantities, the accuracy of gene-finding algorithm in the training and test sets is calculated, which reflects the self-consistency and extrapolating effectiveness of the algorithm. The division of 2958 ORFs into two parts (1958 and 1000) is random. Repeating the above random division procedure three times, we have performed three resubstitution and cross-validation tests. In each case, the constant $c$ is determined by making the false positive rate and false negative rate identical in the resubstitution test. The results of the cross-validation tests are listed in Table 1, where the accuracy is defined as the average of the sensitivity and the specificity. As can be seen from Table 1, the accuracy in each cross-validation test is always greater than 93%. This accuracy is comparable to that obtained by Quentin et al. (1999), based on 745 sequences. Their algorithm performed slightly better than our method, however, so far as we know,

no tests on recognizing all known yeast genes are reported by them.

It should be pointed out that among the 6000 randomly selected intergenic sequences, some sequences may be coding regions containing introns, or code for small proteins or peptides. To examine their influence on the gene-finding result, we shuffle each of the 6000 intergenic sequence 20 000 times to destroy its possible coding structure (yet the overall base composition is the same as the original sequence). Then the shuffled sequences are used as the new negative samples in the algorithm. It is found that the average vector for the negative samples, the constant $c$, as well as the predictive results are very similar to those when the original negative samples are used. This means that the few possible short coding sequences among the negative samples do not affect the predictive result.

There are 223 intron-containing genes of the 1st class in the MIPS database, whose introns have been removed in advance. These ORFs are used as an independent test set to perform another three-fold cross-validation tests. Using the average vectors $\bar{\mathbf{U}}^1$, $\bar{\mathbf{U}}^2$ and the parameter $c$ obtained for each of the three training sets discussed above, the recognition for the 223 sequences is performed. Consequently, the accuracy (defined as the sensitivity in this case) in each test is always greater than 93%, based on the parameters derived from the above three training tests. In other words, a total of six cross-validation tests confirm that the accuracy of the algorithm presented is better than 93%.

### 3.3. Apply the algorithm to find genes in the ORFs of the 2nd–6th classes

After performing the resubstitution and cross-validation tests, the 1958 and 1000 positive samples (true genes) are then merged. The 2958 negative samples are selected randomly from the 6000 intergenic sequences mentioned above. These 2958 positive and 2958 negative samples form a new training set. The random selection is repeated three times. Consequently, we have three combinations. For each combination the positive samples are identical, whereas the negative samples are different each time. The average vectors $\bar{\mathbf{U}}^1$, $\bar{\mathbf{U}}^2$, and the parameter $c$ obtained in each combination are averaged over the three combinations, we find

$$\bar{\mathbf{U}}^1 = (0.2565 \quad -0.0182 \quad 0.0910 \quad -0.0038$$
$$0.1553 \quad 0.2644 \quad -0.0438 \quad -0.0259 \quad 0.2184),$$
$$(9)$$

$$\bar{\mathbf{U}}^2 = (0.0144 \quad 0.0142 \quad 0.2768 \quad -0.0139$$
$$-0.0120 \quad 0.2824 \quad 0.0078 \quad -0.0150 \quad 0.2605),$$
$$(10)$$

$$c = -0.017. \tag{11}$$

We then apply the average vectors $\bar{\mathbf{U}}^1$, $\bar{\mathbf{U}}^2$, and the parameter $c$ listed in Eqs. (9)–(11) to recognizing genes in the ORFs of the 2nd–6th classes in the MIPS database. For each ORF calculate the vector $\mathbf{U} = (u_1, u_2, ..., u_9)^{\mathrm{T}}$, where $u_1$–$u_9$ are defined in Eqs. (1) and (2). Based on the vectors $\mathbf{U}$, $\bar{\mathbf{U}}^1$, $\bar{\mathbf{U}}^2$, and the parameter $c$, calculate the codingness index $\Delta$ using Eq. (6). If $\Delta > 0$, the query ORF is recognized as a coding gene, otherwise, if $\Delta < 0$, the ORF or DNA sequence is recognized as a non-coding one.

It should be pointed out that using the algorithm and the parameters derived from the 1st class ORFs to find genes in the 2nd–6th classes is based on an assumption that both DNA sequences have similar statistical behaviors. This might not be so in some special cases, for example, for some low-expressed genes. In this case, the results of gene-finding in the 2nd–6th class ORFs should be referred to with caution. We hope to see to what extent the assumption is valid, based on a comparison between the results presented here and some future related experiments.

According to the MIPS database, there are 248, 869, 788, 805 and 439 nuclear ORFs of the 2nd–6th classes in the yeast genome. Consequently, 28, 112, 157, 215 and 355 are recognized as non-coding ORFs. The four quantities TP, TN, FP and FN mentioned above can be calculated, based on the above results and the sensitivity and specificity obtained. Compute TP, TN, FP and FN for the 2nd class ORFs in the MIPS database first. The total number of the 2nd class ORFs is 248, in which 28 are recognized as non-coding. Assume that both the sensitivity and specificity are all equal to 93%. We have a set of four linear equations as follows: TP/(TP + FN) = 0.93; TN/(TN + FP) = 0.93; TN + FN = 28 and TP + TN + FP + FN = 248. Solving the above set of equations, we find TP ≈ 219; TN ≈ 12; FP ≈ 1 and FN ≈ 16. The number of real coding ORFs of the 2nd class should be equal to TP + FN ≈ 235. Of the 28 ORFs recognized as non-coding, statistically, 16 (FN) are actually coding. Similar calculations for the 3rd–6th class ORFs are performed. The results are listed in Table 2.

Based on the above results, we re-estimate the number of protein coding genes in the 16 yeast chromosomes. The total number should be equal to: the number of intronless genes in the 1st class (2958) + the number of intron-containing genes in the 1st class (223) + the number of those in the 2nd–6th classes, including intronless and intron-containing genes, recognized by the present algorithm (235 + 810 + 670 + 620 + 63 = 2398, see Table 2). The sum is 5579. Note that the accuracy is actually greater than 93%, so, this

Table 2
The numbers of predicted coding and non-coding ORFs of the 2nd–6th classes

| Class | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|
| Total number of ORFs ($n$)[a] | 248 | 869 | 788 | 805 | 439 |
| TP | 219 | 753 | 623 | 577 | 59 |
| TN | 12 | 55 | 110 | 172 | 350 |
| FP | 1 | 4 | 8 | 13 | 26 |
| FN | 16 | 57 | 47 | 43 | 4 |
| TP+FN[b] | 235 | 810 | 670 | 620 | 63 |
| TN+FN[b] | 28 | 112 | 157 | 215 | 354 |
| TN+FP[b] | 13 | 59 | 118 | 185 | 376 |
| (TN+FP)/$n$[c] | 13/248 = 5.2% | 59/869 = 6.8% | 118/788 = 15.0% | 185/805 = 23.0% | 376/439 = 85.6% |

[a] The mitochondrial ORFs are not included.

[b] TP+FN, TN+FN and TN+FP indicate the numbers of real coding, predicted non-coding and real non-coding ORFs, respectively, based on the assumption that both the sensitivity and specificity of the gene-finding algorithm are equal to 93.0%.

[c] The percentage of the real non-coding ORFs over the total ORFs in this class.

Table 3
The 28 ORFs of the 2nd class in the MIPS database, which are recognized as non-coding[a]

| | | | | | |
|---|---|---|---|---|---|
| YAL004w | YCL069w | YEL004w | YKL008c | YLR034c | YMR118c |
| YAR061w | YDR033w | YER039c | YKL033w-a | YLR046c | YMR279c |
| YBL009w | YDR107c | YER185w | YKR027w | YLR164w | YNL320w |
| YBR161w | YDR276c | YGL054c | YKR105c | YLR176c | |
| YBR210w | YDR384c | YGR131w | YLL051c | YMR040w | |

[a] Of the 28 ORFs listed, statistically, 16 are actually coding. Unfortunately, we cannot identify them at present due to the limited recognition accuracy achieved.

figure should be considered as an up-limit of gene number in the yeast genome. The above estimate of protein coding genes in the yeast genome is about 3.8–7.0% less than that of 5800–6000, which is widely accepted (Goffeau et al., 1996; Winzeler and Davis, 1997; Mewes et al., 1997). The above estimate is based on error analysis, i.e. we have considered the false positive and false negative events in the prediction for each class. So, it should be statistically reliable. As we can see in Table 2, the ratio of non-coding ORFs in each of the 2nd–6th classes is different, which is about 5.2% in 2nd class and 85.6% in 6th class. Meanwhile our estimate (5579) is about 18.7% larger than 4700, estimated recently by another group (Mackiewicz et al., 1999).

The names of the non-coding ORFs of the 2nd–6th classes recognized by the present algorithm are listed in Tables 3–7, respectively. Some of them are actually coding ORFs, but we cannot identify them at present due to the limited accuracy achieved here. We list only the number of such coding ORFs (i.e. FN) in the footnotes of Tables 3–7, respectively, for the 2nd–6th class ORFs in the MIPS database.

### 3.4. Graphic analysis of base composition at different codon positions

As described in Section 2.2, the base composition of an ORF at each codon position can be represented by a point with coordinate ($x_i$, $y_i$, $z_i$) in a three-dimensional space $V_i$, where $i = 1, 2, 3$. Therefore, a set of ORFs is associated with three sets of three-dimensional mapping points. The distribution pattern of these points can be studied in a graphic approach (Zhang and Zhang, 1991; Zhang and Chou, 1994). Using the graphic method, we would like to get some insight into the working mechanism of the algorithm and find the reason why a considerable part of the ORFs in the 2nd–6th classes are actually not coding genes.

As mentioned above, the number of known genes is 3181, and the total number of the ORFs in the 2nd–6th classes is equal to $248 + 869 + 789 + 805 + 447$, i.e. 3158. For comparison, the 6000 intergenic sequences are used as the negative samples. The graphs (Zhang and Zhang, 1991; Zhang and Chou, 1994) corresponding to these data sets are drawn and compared. To save printing space, only the projections onto the $x–y$ or $x–z$ planes are shown here. Consider the known genes

Table 4
The 112 ORFs of the 3rd class in the MIPS database, which are recognized as non-coding[a]

| | | | | | |
|---|---|---|---|---|---|
| YBL089w | YDR249c | YGL160w | YIL166c | YLR064w | YNR056c |
| YBL091c-a | YDR302w | YGL186c | YJL091c | YLR184w | YNR059w |
| YBR074w | YDR303c | YGR023w | YJL170c | YLR251w | YNR063w |
| YBR180w | YDR307w | YGR065c | YJL193w | YLR266c | YOL079w |
| YBR220c | YDR366c | YGR067c | YJR036c | YLR283w | YOL107w |
| YBR293w | YDR387c | YGR077c | YJR124c | YLR311c | YOL119c |
| YCL001w-a | YDR411c | YGR101w | YJR136c | YLR365w | YOL137w |
| YCR023c | YDR413c | YGR284c | YKL037w | YLR394w | YOL152w |
| YCR062w | YEL045c | YHL035c | YKL174c | YML023c | YOL163w |
| YCR087c-a | YEL064c | YHR002w | YKL221w | YMR088c | YOR049c |
| YDL015c | YER048w-a | YHR035w | YKL222c | YMR221c | YOR053w |
| YDL119c | YER097w | YHR048w | YKR030w | YMR245w | YOR292c |
| YDL199c | YER113c | YHR130c | YKR088c | YMR306w | YOR297c |
| YDL206w | YER119c | YHR142w | YKR103w | YNL065w | YOR350c |
| YDL228c | YER184c | YHR181w | YLL005c | YNL109w | YPL125w |
| YDR100w | YFL027c | YIL025c | YLL037w | YNL176c | YPL244c |
| YDR115w | YFL040w | YIL040w | YLL054c | YNL203c | YPR094w |
| YDR119w | YFR057w | YIL054w | YLR010c | YNL275w | |
| YDR205w | YGL104c | YIL088c | YLR050c | YNL305c | |

[a] Of the 112 ORFs listed, statistically, 57 are actually coding. Unfortunately, we cannot identify them at present due to the limited recognition accuracy achieved.

Table 5
The 157 ORFs of the 4th class in the MIPS database, which are recognized as non-coding[a]

| | | | | | |
|---|---|---|---|---|---|
| YAL018c | YDR306c | YGR071c | YJL108c | YML047c | YOL162w |
| YAL034c | YDR438w | YGR110w | YJL147c | YML132w | YOR044w |
| YAR060c | YDR459c | YGR125w | YJR013w | YMR010w | YOR147w |
| YAR068w | YDR492w | YGR212w | YJR044c | YMR034c | YOR175c |
| YBL108w | YDR504c | YGR293c | YJR116w | YMR101c | YOR193w |
| YBL109w | YDR524c | YGR295c | YJR161c | YMR119w | YOR228c |
| YBR004c | YDR525w-a | YHL041w | YJR162c | YMR155w | YOR245c |
| YBR099c | YDR543c | YHL042w | YKL034w | YMR253c | YOR365c |
| YBR147w | YDR544c | YHL044w | YKL219w | YMR324c | YOR390w |
| YBR168w | YEL033w | YHL045w | YKL223w | YMR326c | YPL087w |
| YBR183w | YEL067c | YHL048w | YKL225w | YNL008c | YPL165c |
| YBR300c | YER072w | YHR054c | YKR051w | YNL026w | YPL189w |
| YBR302c | YER188c-a | YHR133c | YKR106w | YNL101w | YPL229w |
| YCL002c | YFL015c | YHR162w | YLL023c | YNL156c | YPL246c |
| YCL038c | YFL062w | YHR212c | YLL031c | YNL297c | YPL257w |
| YCL073c | YFL063w | YHR214w-a | YLR023c | YNL326c | YPL264c |
| YCR102w-a | YFL065c | YHR217c | YLR036c | YNL336w | YPL279c |
| YCR103c | YFL068w | YIL029c | YLR047c | YNL337w | YPR071w |
| YDL123w | YFR012w | YIL089w | YLR156c | YNR062c | YPR114w |
| YDL183c | YGL010w | YIL090w | YLR159c | YNR075c | YDR367w* |
| YDL248w | YGL041c | YIL174w | YLR161w | YNR077c | YMR292w* |
| YDR018c | YGL084c | YIL175w | YLR241w | YOL002c | YOL047c* |
| YDR066c | YGL124c | YIR040c | YLR246w | YOL003c | |
| YDR084c | YGL260w | YIR043c | YLR414c | YOL048c | |
| YDR105c | YGL263w | YIR044c | YLR463c | YOL092w | |
| YDR126w | YGR015c | YJL062w | YML033w | YOL101c | |
| YDR131c | YGR016w | YJL097w | YML036w | YOL129w | |

[a] Of the 157 ORFs listed above, 154 are intronless and three are intron-containing (marked with *). Note that of the 157 ORFs listed, statistically, 47 are actually coding. Unfortunately, we cannot identify them at present due to the limited recognition accuracy achieved.

Table 6
The 215 ORFs of the 5th class in the MIPS database, which are recognized as non-coding[a]

| | | | | | |
|---|---|---|---|---|---|
| YAL008w | YDL231c | YGL057c | YJR041c | YLR404w | YNL311c |
| YAL064w | YDR015c | YGL138c | YJR120w | YML003w | YNL324w |
| YAL066w | YDR024w | YGL188c | YJR157w | YML038c | YOL024w |
| YAR030c | YDR029w | YGL230c | YKL044w | YML084w | YOL038c-a |
| YAR040c | YDR042c | YGR026w | YKL051w | YML090w | YOL072w |
| YAR047c | YDR065w | YGR149w | YKL097c | YML107c | YOL118c |
| YAR053w | YDR102c | YGR168c | YKL102c | YML122c | YOL160w |
| YAR064w | YDR141c | YGR226c | YKL158w | YMR003w | YOL166c |
| YAR069c | YDR179w-a | YGR290w | YKL162c | YMR007w | YOR015w |
| YAR070c | YDR215c | YGR291c | YKR032w | YMR057c | YOR024w |
| YBL048w | YDR274c | YHL005c | YKR065c | YMR071c | YOR029w |
| YBL049w | YDR278c | YHL037c | YKR103c | YMR082c | YOR068c |
| YBL071c | YDR315c | YHR067w | YLL007c | YMR103c | YOR072w |
| YBR013c | YDR319c | YHR078w | YLL014w | YMR122c | YOR080w |
| YBR027c | YDR344c | YHR095w | YLL030c | YMR141c | YOR152c |
| YBR058c-a | YDR350c | YHR139c-a | YLL033w | YMR148w | YOR183w |
| YBR085c-a | YDR396w | YHR140w | YLL042c | YMR151w | YOR268c |
| YBR096w | YDR437w | YHR173c | YLL059c | YMR163c | YOR314w |
| YBR126w-a | YDR524w-a | YIL012w | YLR042c | YMR187c | YOR364w |
| YBR144c | YDR525w | YIL028w | YLR049c | YMR191w | YOR376w |
| YBR292c | YEL010w | YIL037c | YLR111w | YMR252c | YOR392w |
| YCL056c | YEL014c | YIL071c | YLR112w | YMR254c | YPL041c |
| YCL057c-a | YEL059w | YIL086c | YLR122c | YMR258c | YPL052w |
| YCL058c | YER044c | YIL152w | YLR124w | YMR259c | YPL056c |
| YCR001w | YER050c | YIR020c | YLR199c | YMR320w | YPL066w |
| YCR006c | YER066c-a | YIR020c-a | YLR255c | YNL017c | YPL162c |
| YCR022c | YER091c-a | YJL027c | YLR264c-a | YNL038w | YPL200w |
| YCR025c | YER135c | YJL028w | YLR265c | YNL122c | YPR012w |
| YCR043c | YER140w | YJL064w | YLR267w | YNL143c | YPR014c |
| YCR085w | YER172c-a | YJL077c | YLR296w | YNL146w | YPR064w |
| YDL027c | YFL019c | YJL118w | YLR312c | YNL150w | YPR098c |
| YDL054c | YFL021c-a | YJL136w-a | YLR366w | YNL174w | YPR153w |
| YDL089w | YFR035c | YJL163c | YLR376c | YNL179c | YPR170c |
| YDL162c | YFR042w | YJL215c | YLR381w | YNL211c | YPR170w-a |
| YDL180w | YFR054c | YJR011c | YLR400w | YNL269w | YDR535c* |
| YDL196w | YGL006w-a | YJR023c | YLR402w | YNL303w | |

[a] Of the 215 ORFs listed above, 214 are intronless and one is intron-containing (marked with *). Note that of the 215 ORFs listed, statistically, 43 are actually coding. Unfortunately, we cannot identify them at present due to the limited recognition accuracy achieved.

first. The projections of the 3181 mapping points onto the $x$–$y$ planes for the base composition at the 1st and 2nd codon positions are shown in Fig. 1(a) and (b), respectively. The projection of these mapping points onto the $x$–$z$ plane for the base composition at the 3rd codon position is shown in Fig. 1(c). As can be seen, the distribution patterns of the projection points of the three plots are considerably different. In Fig. 1(a), most of the points are situated at the region where $a_1 > t_1$ and $g_1 > c_1$, i.e. purine bases are predominant at the first codon position. In Fig. 1(b), most of the points are situated at the region lacking G. In Fig. 1(c), most of the points are at the region where A and T are predominant. In summary, the preference codons of the yeast protein coding genes are of the $R\bar{G}W$ type, where R, $\bar{G}$ and W are the bases of purine, non-G and A/T, respec-

tively. For comparison, the projection of the mapping points onto the $x$–$y$ plane for the base composition at the 1st 'codon' position of the 6000 intergenic sequences is shown in Fig. 2. Note that the 'codon' in an intergenic sequence is meaningless. Because all the intergenic sequences selected here begin with ATG, for example, ATGGCGCAT…, the bases A, G, C… are defined as at the first 'codon' position and so forth. Since the distribution patterns of the points at the 2nd and 3rd 'codon' positions are almost identical with that at the 1st 'codon' position, they are not shown here. Therefore, the 'codons' of the intergenic sequences are of the type NNN, where N indicates any base.

It has been suggested that the first, second and third position of the codons are associated, respectively, with the biosynthetic pathway, hydrophobicity pattern, and

Table 7
The 355 ORFs of the 6th class in the MIPS database, which are recognized as non-coding[a]

| | | | | | |
|---|---|---|---|---|---|
| YAL034c-b | YDR114c | YGL182c | YJL135w | YLR428c | YOR102w |
| YAL042c-a | YDR133c | YGL193c | YJL142c | YLR434c | YOR105w |
| YAL056c-a | YDR136c | YGL204c | YJL150w | YLR444c | YOR121c |
| YBL012c | YDR149c | YGL214w | YJL152w | YLR458w | YOR135c |
| YBL053w | YDR154c | YGL217c | YJL169w | YLR465c | YOR146w |
| YBL062w | YDR157w | YGL218w | YJL175w | YML009w-a | YOR169c |
| YBL065w | YDR187c | YGL239c | YJL182c | YML012c-a | YOR170w |
| YBL070c | YDR199w | YGR011w | YJL188c | YML031c-a | YOR199w |
| YBL073w | YDR203w | YGR018c | YJL202c | YML034c-a | YOR200w |
| YBL077w | YDR220c | YGR025w | YJL220w | YML047w-a | YOR203w |
| YBL083c | YDR230w | YGR039w | YJR018w | YML057c-a | YOR218c |
| YBL094c | YDR241w | YGR045c | YJR020w | YML089c | YOR225w |
| YBL107w-a | YDR269c | YGR050c | YJR037w | YML094c-a | YOR235w |
| YBR051w | YDR271c | YGR051c | YJR038c | YML099w-a | YOR248w |
| YBR064w | YDR290w | YGR064w | YJR071w | YML116w-a | YOR263c |
| YBR089w | YDR340w | YGR069w | YJR087w | YMR046w-a | YOR277c |
| YBR109w-a | YDR355c | YGR073w | YJR128c | YMR052c-a | YOR282w |
| YBR113w | YDR360w | YGR107w | YKL030w | YMR075c-a | YOR300w |
| YBR116c | YDR401w | YGR114c | YKL036c | YMR086c-a | YOR309c |
| YBR124w | YDR417c | YGR115c | YKL053w | YMR119w-a | YOR325w |
| YBR134w | YDR426c | YGR122c-a | YKL076c | YMR135w-a | YOR331c |
| YBR178w | YDR431w | YGR137w | YKL083w | YMR153c-a | YOR333c |
| YBR206w | YDR442w | YGR139w | YKL111c | YMR158c-b | YOR345c |
| YBR224w | YDR445c | YGR151c | YKL115c | YMR158w-a | YOR379c |
| YBR226c | YDR455c | YGR164w | YKL118w | YMR172c-a | YPL034w |
| YBR232c | YDR467c | YGR176w | YKL123w | YMR193c-a | YPL035c |
| YBR266c | YDR491c | YGR182c | YKL131w | YMR290w-a | YPL044c |
| YBR277c | YDR509w | YGR190w | YKL136w | YMR304c-a | YPL073c |
| YCL006c | YDR521w | YGR219w | YKL147c | YMR306c-a | YPL102c |
| YCL023c | YDR526c | YGR228w | YKL153w | YMR316c-a | YPL114w |
| YCL041c | YEL075w-a | YGR242w | YKL162c-a | YNL013c | YPL136w |
| YCL065w | YER006c-a | YGR259c | YKL169c | YNL028w | YPL182c |
| YCR018c-a | YER046w-a | YGR265w | YKL177w | YNL089c | YPL185w |
| YCR041w | YER067c-a | YHL002c-a | YKL202w | YNL105w | YPL205c |
| YCR064c | YER084w | YHL006w-a | YKR033c | YNL114c | YPL238c |
| YCR087w | YER138w-a | YHL030w-a | YKR040c | YNL120c | YPL261c |
| YDL009c | YER145c-a | YHL046w-a | YKR047w | YNL170w | YPR002c-a |
| YDL011c | YER148w-a | YHR049c-a | YLL020c | YNL171c | YPR038w |
| YDL016c | YER165c-a | YHR056w-a | YLR101c | YNL184c | YPR039w |
| YDL023c | YER181c | YHR063w-a | YLR123c | YNL198c | YPR050c |
| YDL026w | YFL012w-a | YHR125w | YLR140w | YNL205w | YPR053c |
| YDL032w | YFL013w-a | YHR145c | YLR169w | YNL226w | YPR077c |
| YDL034w | YFL032w | YIL060w | YLR171w | YNL228w | YPR087w |
| YDL041w | YFR036w-a | YIL066w-a | YLR198c | YNL235c | YPR092w |
| YDL050c | YFR056c | YIL068w-a | YLR217w | YNL266w | YPR099c |
| YDL062w | YGL024w | YIL071w-a | YLR230w | YNL276c | YPR126c |
| YDL068w | YGL042c | YIL100c-a | YLR232w | YNL296w | YPR130c |
| YDL071c | YGL052w | YIL156w-a | YLR252w | YNL319w | YPR136c |
| YDL094c | YGL072c | YIL163c | YLR261c | YNR005c | YPR142c |
| YDL151c | YGL074c | YIL171w-a | YLR269c | YNR025c | YPR146c |
| YDL152w | YGL088w | YIR017w-a | YLR279w | YOL013w-a | YPR150w |
| YDL158c | YGL102c | YIR023c-a | YLR282c | YOL013w-a | YPR177c |
| YDL172c | YGL109w | YJL009w | YLR294c | YOL035c | YBR090c* |
| YDL187c | YGL118c | YJL015c | YLR302c | YOL037c | YER014c-a* |
| YDL221w | YGL132w | YJL022w | YLR317w | YOL099c | YLR202c* |
| YDR008c | YGL149w | YJL032w | YLR322w | YOL106w | |
| YDR034c-a | YGL152c | YJL067w | YLR334c | YOL134c | |
| YDR048c | YGL165c | YJL075c | YLR339c | YOL150c | |
| YDR053w | YGL168w | YJL086c | YLR358c | YOR041c | |
| YDR112w | YGL177w | YJL120w | YLR379w | YOR082c | |

[a] Of the 355 ORFs listed above, 352 are intronless and three are intron-containing (marked with *). Note that of the 355 ORFs listed, statistically, four are actually coding. Unfortunately, we cannot identify them at present due to the limited recognition accuracy achieved.
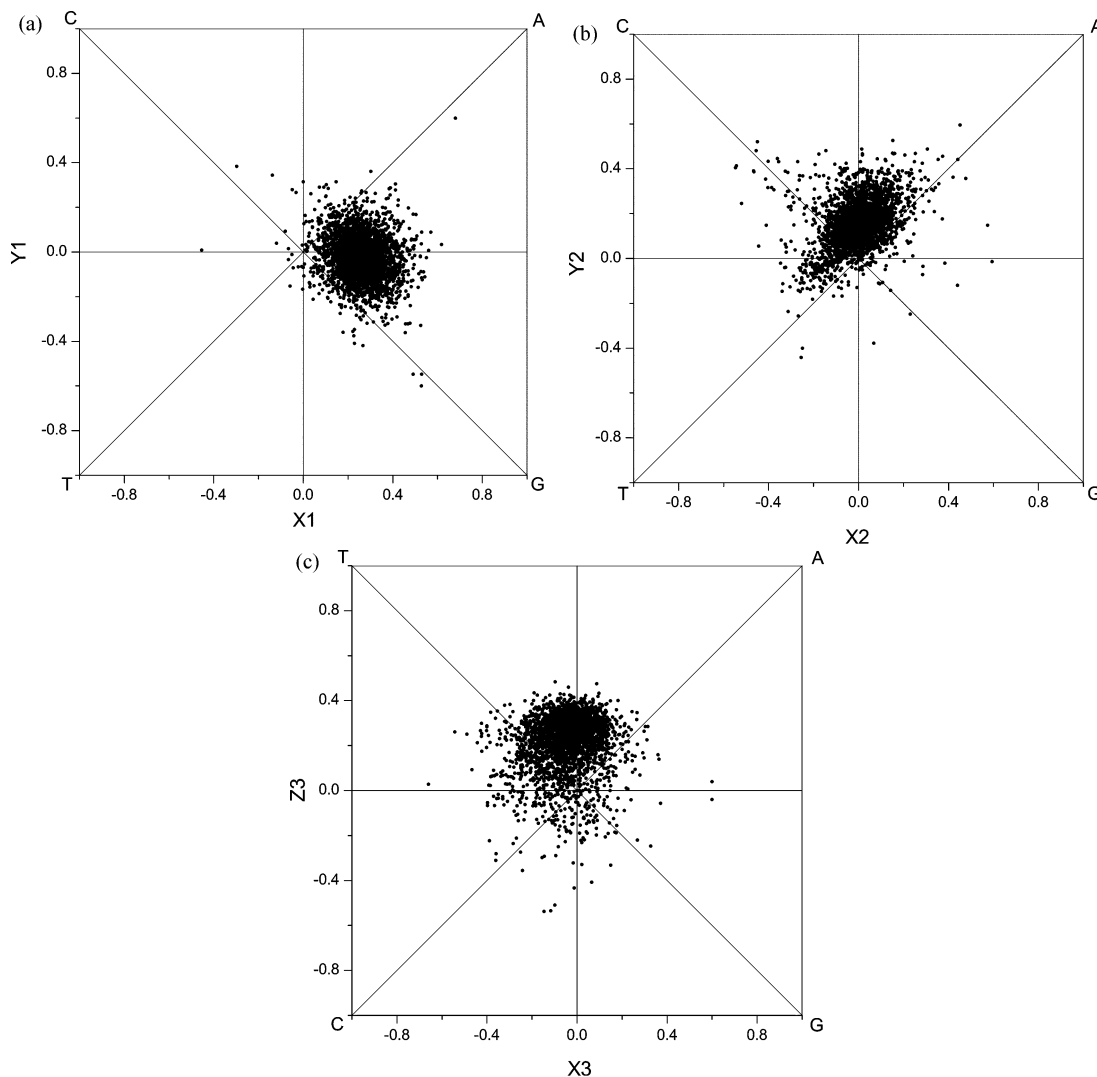
Fig. 1. Distributions of the mapping points for base composition at three codon positions of the 3181 protein coding genes in the yeast genome in a three-dimensional space. (a) Projection onto the $x$–$y$ plane for the 1st codon position; (b) projection onto the $x$–$y$ plane for the 2nd codon position; and (c) projection onto the $x$–$z$ plane for the 3rd codon position. For more detailed explanation regarding the graph, refer to Zhang and Zhang (1991) or Zhang and Chou (1994).

the alpha helix or beta strand forming potentiality of the coded amino acids (Siemion and Siemion, 1994; Taylor and Coates, 1989). Recently, it is reported that there is strong correlation between the base frequencies in the second codon position of genes and the corresponding secondary structures in the encoded proteins (Gupta et al., 2000; Chiusano et al., 2000). Chiusano et al. (2000) attributed this relation to the hydrophobic and hydrophilic amino acids encoded by codons having U or A, respectively, in their second codon site. It is also supposed that the specific codon choice is functionally needed in mRNA–rRNA interaction in ribosome, which is responsible for monitoring the correct reading frame during translation (Trifonov, 1987; Lagunez-

Otero and Trifonov, 1992; Lobry and Gautier, 1994). These mean that the three codon positions are associated with different biological functions and the base choices at these positions are usually specific. In the case of yeast genome, the preferred codon usage pattern is RḠW, and the difference of the two codon types, i.e. RḠW and NNN, forms the basis to distinguish between coding and non-coding sequences. The present algorithm is based on the difference of mapping point distribution patterns between the two kinds of sequences. As we can see, the present algorithm works well.

For comparison, the projection of the mapping points for the base composition at the 1st codon posi-

tion of the 3158 ORFs of the 2nd–6th classes in the MIPS database, is shown in Fig. 3. Obviously, the distribution pattern of Fig. 3 is in between those of Fig. 1(a) and Fig. 2, indicating that a part of ORFs of the 2nd–6th classes in the MIPS database are actually of non-coding ones. This fact implies that the codons in some ORFs of the 2nd–6th classes in the MIPS database are of the RGW type, whereas the 'codons' in some other ORFs are of the NNN type. The latter 'codons' do not code for any realistic proteins. This is the reason why 23% (refer to the last row in Table 2) of the ORFs of the 2nd–6th classes in the MIPS database are non-coding.

In fact, the division between coding and non-coding regions can be seen in a more intuitive manner by a principal component analysis (PCA). PCA defines a rotation of the variables of given data. The first derived direction (a linear combination of the variables) is chosen to maximize the standard deviation of the derived variable, the second to maximize the standard deviation among directions un-correlated with the first, and so forth. For details about this method refer to Dillon and Goldstein (1984). For the data set comprising the values of $u_1–u_9$ of 3181 known genes and 6000 intergenic sequences, a plot based on the two most important axes using the PCA is shown in Fig. 4. The first and the second axis account for 34.3 and 15.9% of



Fig. 3. Distributions of the mapping points for base composition of the 3158 ORFs of the 2nd–6th classes in the MIPS database in a three-dimensional space. Projection onto the $x–y$ plane for the 1st codon position. Specially note that the distribution pattern here is in between those in Fig. 1(a) and Fig. 2, indicating that a part of the ORFs of the 2nd–6th classes in the MIPS database is actually non-coding. See the legend of Fig. 1 for explanation of the graph.

the total inertia of the nine-dimensional space, respectively, and no other axis accounts for more than 10%. The variation in the second axis is mainly because of a small number of outliers, which are mostly short genes or ORFs. The two principal axes are responsible for separating the coding and non-coding sequences into two clusters. The coding sequences are represented by close cycles and the non-coding sequences are represented by open cycles. As we can see, the two clusters appear to be distinct, with quite little overlap. The closeness of any two regions in Fig. 4 reflects the similarities of their base frequencies at the three codon positions, implying that the base choices at the three codon positions for most genes are quite different from those of non-coding sequences.

## 4. Conclusions

A simple gene-finding algorithm with high accuracy (93%) for the yeast genome is presented in this paper. Six-fold cross-validation tests confirm the above accuracy. Using the algorithm, it is found that 751 ORFs (about 23.8%) of the 2nd–6th classes classified in the MIPS database are likely non-coding. The total number of protein coding genes in the 16 yeast chromosomes is estimated to be less than or equal to 5579. This estimate is based on the assumption that the DNA sequences coding for proteins in the 1st class ORFs have similar
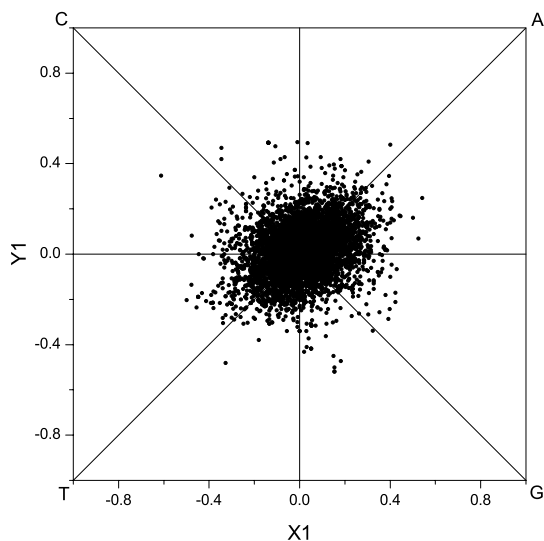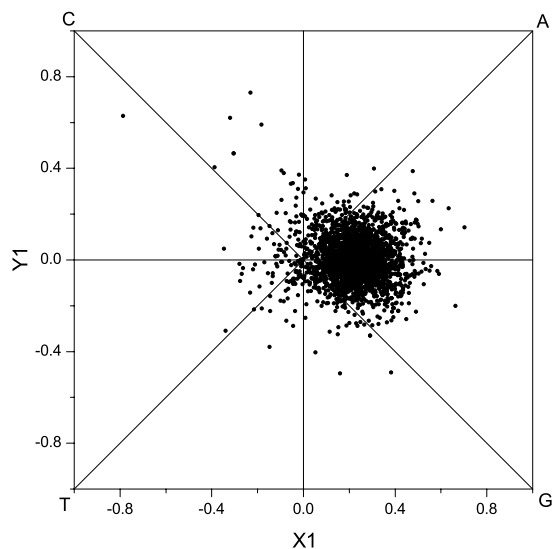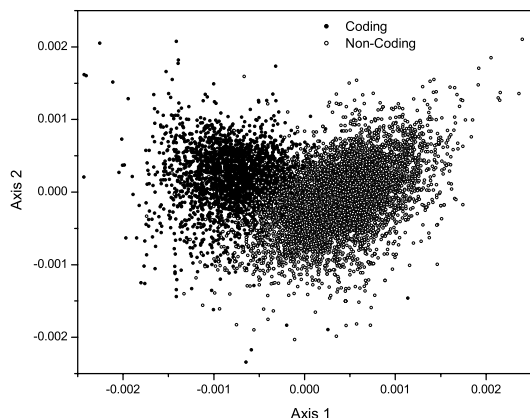


Fig. 2. Distributions of the mapping points for base composition of about 6000 intergenic sequences in the yeast genome in a three-dimensional space. Projection onto the $x–y$ plane for the 1st 'codon' position. It should be pointed out that the distribution patterns for the 2nd and 3rd 'codon' positions (not shown here) are almost identical to that for the 1st 'codon' position. Note that the 'codon' in an intergenic sequence is meaningless. The definition of 'codon' positions in an intergenic sequence is arbitrary. See the legend of Fig. 1 for explanation of the graph.

Fig. 4. The distribution of points based on the two most important axes using the principal component analysis of the nine variables $u_1$–$u_9$ for the 3181 known nuclear genes and 6000 intergenic sequences. Axis 1 is a derived direction (a linear combination of the variables $u_1$–$u_9$) chosen to maximize the standard deviation of the derived variable, and Axis 2 is the direction to maximize the standard deviation among directions un-correlated with the first. For the values of $u_1$–$u_9$ of 3181 known genes and 6000 intergenic sequences, Axes 1 and 2 account for 34.3 and 15.9% of the total inertia of the nine-dimensional space, respectively. The close cycles indicate the known genes while the open cycles represent the non-coding sequences. Note that the two clusters appear to be distinct, with quite little overlap. This fact constitutes the basis of the present algorithm.

statistical properties to those coding for proteins in the 2nd–6th class ORFs. We hope to see to what extent the above assumption is valid, relying on a comparison between the results presented here and some future related experiments. The working mechanism of the present algorithm is studied in detail by a graphic approach. It is found that the preference codons of the yeast protein coding genes are of the $R\bar{G}W$ type, where R, $\bar{G}$ and W are the bases of purine, non-G and A/T, respectively, whereas the 'codons' of the intergenic sequences are of the type NNN, where N indicates any base. The difference of the two codon types, i.e. $R\bar{G}W$ and NNN, forms the basis of the algorithm to distinguish between coding and non-coding sequences. Mathematically, the present algorithm is based on the difference of mapping point distributions between coding and non-coding sequences. The algorithm can be extended to find genes in any prokaryotic genome, but it cannot be directly applied to higher eukaryotic genomes with more introns.

## Acknowledgements

## References

Benetzen, J., Benjamin, D.H., 1982. Codon selection in yeast. J. Biol. Chem. 257, 3026–3031.

Borodovsky, M.Y., McIninch, J.D., 1993. GeneMark: parallel gene recognition for both DNA strands. Comput. Chem. 17, 123–153.

Bowman, S., Churcher, C., Badcock, K., Brown, D., Chillingworth, T., Connor, R., Dedman, K., Gentels, S., Hamlin, N., Hunt, S., et al., 1997. The nucleotide sequence of *Saccharomyces cerevisiae* chromosome XIII. Nature. (Suppl.) 387, 90–93.

Burset, M., Guigo, R., 1996. Evaluation of gene structure prediction programs. Genomics 34, 353–367.

Bussey, H., Kabak, D.B., Zhong, W., Vo, D.T., Cloak, M.W., Fortin, N., Hall, J., Ouellette, B.F., Keng, T., Barton, A.B., et al., 1995. The nucleotide sequence of chromosome I from *Saccharomyces cerevisiae*. Proc. Natl. Acad. Sci. 92, 3809–3813.

Bussey, H., Storms, R.K., Ahmed, A., Albermann, K., Allen, E., Ansorge, W., Araujo, R., Aparicio, A., Barrell, B., Badcock, K., et al., 1997. The nucleotide sequence of *Saccharomyces cerevisiae* chromosome XVI. Nature. (Suppl.) 387, 103–105.

Chiusano, M.L., Alvarez-Valin, F., Giulio, M.D., D'Onofrio, G., Ammirato, G., Colonna, G., Bernardi, G., 2000. Second codon positions of genes and the secondary structures of proteins. Relationships and implications for the origin of the genetic code. Gene 261, 63–69.

Churcher, C., Bowman, S., Badcock, K., Bankier, A., Brown, D., Chillingworth, T., Connor, R., Delvin, K., Gentles, S., Hamlyn, N., et al., 1997. The nucleotide sequence of *Saccharomyces cerevisiae* chromosome IX. Nature. (Suppl.) 387, 84–87.

Dietrich, F.S., Mulligan, J., Hennessy, K., Yelton, M.A., Allen, E., Araujo, R., Aviles, E., Berno, A., Brennan, T., Carpenter, J., et al., 1997. The nucleotide sequence of *Saccharomyces cerevisiae* chromosome V. Nature. (Suppl.) 387, 78–81.

Dillon, W.R., Goldstein, M., 1984. Multivariate Analysis, Methods and Applications. Wiley, New York.

Dujon, B., Alexandraki, D., Andre, B., Ansorge, W., Baladron, V., Ballesta, J.P.G., Banrevi, A., Bolle, P.A., Bolotin-Fukuhara, M., Bossier, P., et al., 1994. Complete DNA sequence of yeast chromosome XI. Nature. (Suppl.) 369, 371–378.

Dujon, B., Albermann, K., Aldea, M., Alexandraki, D., Ansorge, W., Arino, J., Benes, V., Bohn, C., Bolotin-Fukuhara, M., Bordonne, R., et al., 1997. The nucleotide sequence of *Saccharomyces cerevisiae* chromosome XV. Nature. (Suppl.) 387, 98–102.

Feldmann, H., et al., 1994. Complete DNA sequence of yeast chromosome II. EMBO J. 13, 5793–5809.

Fickett, J.W., 1996. Finding genes by computer: the state of the art. Trends Genet. 12, 316–320.

Galibert, F., Alexandraki, D., Baur, A., Boles, E., Chalwatzis, N., Chuat, J.-C., Coster, F., Cziepluch, C., De Haan, M., Domde, H., et al., 1996. Complete nucleotide sequence of *Saccharomyces cerevisiae* chromosome X. EMBO J. 15, 2031–2049.

Goffeau, A., Barrel, B.G., Bussey, H., Davis, R.W., Dujon, B., Feldmann, H., Galibert, F., Hoheisel, J.D., Jacq, C., Johnston, M., Louis, E.J., Mewes, H.W., Murakami, Y., Philippsen, P., Tettlin, H., Oliver, S.G., 1996. Science 274, 546.

Gupta, S.K., Majumdar, S., Bhattacharya, T.K., Ghosh, T.C., 2000. Studies on the relationships between the synonymous codon usage and protein secondary structural units. Biochem. Biophys. Res. Commun. 269, 692–696.

Jacq, C., Alt-Morbe, J., Andre, B., Arnold, W., Bahr, A., Ballesta, J.P.G., Bargues, M., Baron, L., Becker, A., Biteau, N., et al., 1997. The nucleotide sequence of *Saccharomyces cerevisiae* chromosome IV. Nature. (Suppl.) 387, 75–78.

Johnston, M., Andrews, S., Brinkman, R., Cooper, J., Ding, H., Dover, J., Du, Z., Favello, A., Fulton, L., Gattung, S., et al., 1994. Complete nucleotide sequence of *Saccharomyces cerevisiae* chromosome VIII. Science 265, 2077–2082.

Johnston, M., Hillier, L., Riles, L., Albermann, K., Andre, B., Ansorge, W., Benes, V., Bruckner, M., Delius, H., et al., 1997. The nucleotide sequence of *Saccharomyces cerevisiae* chromosome XII. Nature. (Suppl.) 387, 87–90.

Lagunez-Otero, J., Trifonov, E.N., 1992. mRNA periodical infrastructure complementary to the proof-reading site in the ribosome. J. Biomolec. Struct. Dyn. 10, 455–464.

Lobry, J.R., Gautier, C., 1994. Hydrophobicity, expressivity and aromaticity are the major trends of amino-acid usage in 999 *Escherichia coli* chromosome-encoded genes. Nucleic Acids Res. 22, 3174–3180.

Mackiewicz, P., Kowalczuk, M., Gierlik, A., Dudek, M.R., Cebrat, S., 1999. Origin and properties of non-coding ORFs in the yeast genome. Nucleic Acids Res. 27, 3503–3509.

Mewes, H.W., Albermann, K., Bahr, M., Frishman, D., Gleissner, A., Hani, J., Heumann, K., Kleine, K., Maierl, A., Oliver, S.G., Pfeiffer, F., Zollner, A., 1997. Overview of the yeast genome. Nature. (Suppl.) 387, 7–8.

Mewes, H.W., Heumann, K., Kaps, A., Mayer, K., Pfeiffer, F., Stocker, S., Frishman, D., 1999. MIPS: a database for protein sequences and complete genomes. Nucleic Acids Res. 27, 44–48.

Murakami, Y., Naitou, M., Hagiwara, H., Shibata, T., Ozawa, M., Sasanuma, S.I., Sasanuma, M., Tsuchiya, Y., Soeda, E., Yokoyama, K., et al., 1995. Analysis of the nucleotide sequence of chromosome VI from *Saccharomyces cerevisiae*. Nature Genet. 10, 262–268.

Oliver, S.G., van der Aart, Q.J.M., Agosoni-Carbone, M.L., Aigle, M., Alberghina, L., Alexandraki, D., Antoine, G., Anwar, R., Ballesta, J.P.G., Benit, P., et al., 1992. The complete DNA sequence of yeast chromosome III. Nature 357, 38–46.

Philippsen, P., Kleine, K., Pohlmann, R., Dusterhoft, A., Hamberg, K., Hagemann, J.H., Obermaier, B., Urrestarazu, L.A., Aert, R., Albermann, K., et al., 1997. The nucleotide sequence of *Saccharomyces cerevisiae* chromosome XIV and its evolutionary implications. Nature. (Suppl.) 387, 93–98.

Quentin, Y., Voiblet, C., Martin, F., Fichant, G., 1999. Protein-coding region discovery in organisms underrepresented in databases. Comput. Chem. 23, 209–217.

Sharp, P.M., Li, W.-H., 1987. The codon adaptation index—a measure of directional synonymous codon usage bias, and its potential application. Nucleic Acids Res. 15, 1281–1295.

Siemion, I.Z., Siemion, P.J., 1994. The informational context of the third base in amino acid codons. Biosystems 33, 39–48.

Taylor, F.J., Coates, D., 1989. The code within the codons. Biosystems 22, 177–187.

Tettelin, H., Agostoni-Carbone, M.L., Albermann, K., Albers, M., Arroyo, J., Backes, U., Barreiros, T., Bertani, I., Bjourson, A.J., Brucker, M., et al., 1997. The nucleotide sequence of *Saccharomyces cerevisiae* chromosome VII. Nature. (Suppl.) 387, 81–84.

Trifonov, E.N., 1987. Translation framing code and frame-monitoring mechanism as suggested by the analysis of mRNA and 16 S rRNA nucleotide sequences. J. Mol. Biol. 194, 643–652.

Winzeler, E.A., Davis, R.W., 1997. Functional analysis of the yeast genome. Curr. Opin. Genet. Dev. 7, 771–776.

Zhang, C.T., Chou, K.C., 1994. A graphic approach to analyzing codon usage in 1562 *E. coli* protein coding sequences. J. Mol. Biol. 238, 1–8.

Zhang, C.T., Zhang, R., 1991. Analysis of distribution of bases in the coding sequences by a diagrammatic technique. Nucleic Acids Res. 19, 6313–6317.