

Fig. 6. (A) ChiM93-10 and (B) ChiM93-12 mice simultaneously infected with genotypes A and G.

12) received simultaneous inoculation with 10^5 copies each of HBV/A (A2_JPN strain) and HBV/G (G_US1646 strain). The ChiM93-10 mouse developed HBV/G DNA 17 weeks after inoculation, 9 weeks since HBV/A DNA had increased to $>10^7$ copies/ml (Fig. 6A). HBV/G DNA increased to the level of total HBV DNA at week 21, thereby indicating that by then, HBV/G had taken over HBV/A almost completely.

For the reasons unknown, infection with HBV/G was not established in the ChiM93-12 mouse simultaneously coinfecting with HBV/A (Fig. 6B), although it was infected with HBV/A in levels by some 2 logs lower (10^7 copies/ml) than the ChiM93-10 mouse. Serum levels of human albumin in the ChiM93-12 mouse (mean, 2.1×10^6 ng/ml) were much lower than the other chimeric mice used in this study (mean, 4.7×10^6 ng/ml). Thus, a lower extent of repopulation with human hepatocytes may have prohibited active replication of HBV/A. This would be a prerequisite to infection with HBV/G at high levels.

Coinfection of Mice with HBV/A and HBV/G by Inoculation with a Mouse Passage of G-on-A Superinfection. Three ChiM mice (ChiM169-8, ChiM133-3, and ChiM133-6) received serum from a ChiM92-9 mouse with G-on-A superinfection taken at week 26, when HBV/G had almost replaced HBV/A (Fig. 3A). Profiles of HBV/A and HBV/G, after inoculation with 10^5 copies of HBV DNA, were similar among the mice (Fig. 7A-C). HBV/G DNA was detected at week 1 in levels comparable to those of total HBV DNA. Despite receiving the inoculation with a mouse passage containing HBV/G, in copies by 5 logs greater than those of HBV/A,

HBV/G DNA decreased thereafter and stayed >1 log lower than total HBV DNA until week 7. Since week 4, HBV/G started to increase and replaced HBV/A almost completely until weeks 10-12, and continued to do so through weeks 19-22 of the observation (Fig. 7A).

Cloning and Sequencing HBV DNA in Chimeric Mice Coinfected with HBV/A and HBV/G. HBV DNA clones from sera of ChiM92-9 sampled at 26 weeks (Fig. 4A) and ChiM169-8 inoculated with serum passage in it (Fig. 7A) included those of HBV/A and G invariably. They confirmed the results of real-time detection PCR and PCR-RFLP and did not possess any mutations in comparison with the original inoculum of either genotype. No recombinations between HBV/A and G were detected, either. At least 5 clones of each genotype were propagated and sequenced in both sera.

Cotransfection of Huh7 Cells with Plasmids Carrying the Core Gene of Genotype A and the Entire Genome of Genotype G. Huh7 cells were transfected with 2 plasmids that were pcDNA_core clones that expressed the core protein of genotype A2, under the control of cytomegalovirus promoter, and the pUC19/G clone incorporated with 1.24-fold the genome of genotype G. Transfection only with genotype G induced its replication in a weak level (Fig. 8). When Huh7 cells were cotransfected with the genotype G clone and the genotype A core clone, however, the replication was enhanced in a dose-dependent manner.

Liver Pathology of ChiM Mice Infected with HBV/A and/or HBV/G. Figure 9 shows the histology of liver in representative ChiM mice either simultaneously coinfecting with genotypes A and G (viremia of only ge-

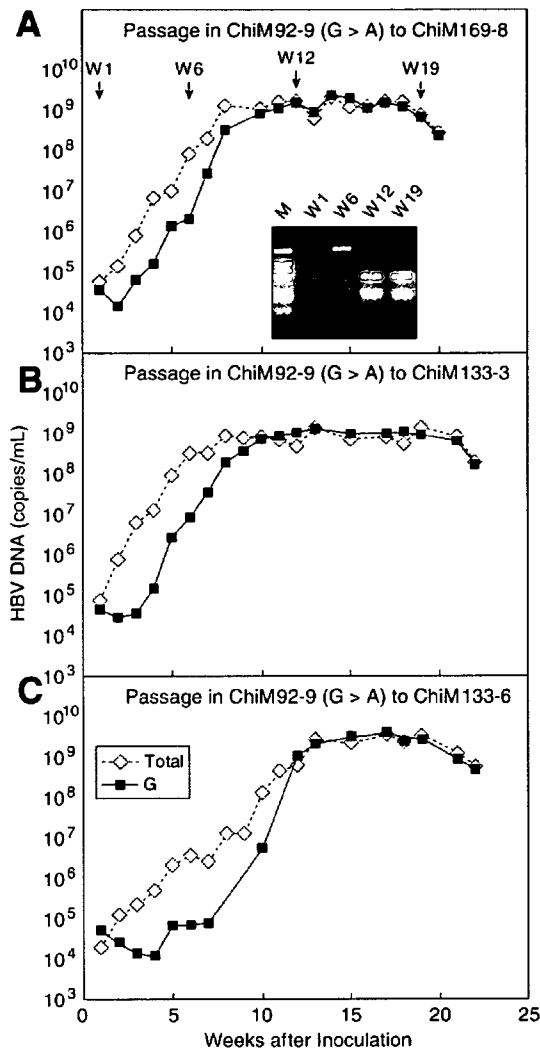


Fig. 7. (A) ChiM169-8, (B) ChiM133-3, and (C) ChiM133-6 mice inoculated with a serum passage from a mouse coinfecting with genotypes A and G (ChiM92-9 in Fig. 4A).

notype A in ChiM93-12) or superinfected with genotypes G-on-A (ChiM92-9) and monoinfected with genotype G (ChiM92-3) during 32-39 weeks. HBV infection was demonstrated by double staining for HBcAg and human albumin (Supplementary Fig. 2). The mouse coinfecting with genotypes A and G revealed steatosis of hepatocytes with hematoxylin-eosin stain and fibrosis of stage 2 (F2) with Masson's trichrome stain. In contrast, the mice monoinfected with genotype A (ChiM93-12) or G (ChiM92-3) had neither steatosis nor fibrosis. Table 1 summarizes the liver pathology of all autopsied mice. Steatosis in 30%-80% of repopulated human hepatocytes and stage F1-F2 fibrosis were observed in the majority of mice superinfected or coinfecting with genotypes G and A or C.

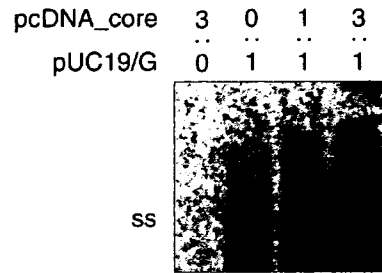


Fig. 8. Trans-complementation of the core gene of genotype A enhanced replication of genotype G. Huh7 cells were cotransfected with plasmids constructed with 1.24-fold the genome of genotype G (pUC19/G) and plasmids expressing the core gene of genotype A (pcDNA_core) in an increasing ratio. Gel strips were Southern-blotted by the complete HBV probe of genotype G. The far left lane represents negative control with pcDNA_core alone. The migration position of single-stranded (ss) HBV DNA is indicated on the left.

Discussion

Using ChiM mice infected with pedigreed HBV DNA in the standardized copy number, we have determined early viral dynamics of HBV/G in detail. Due to constraints on securing ChiM mice with a satisfactory rate of replacement for human hepatocytes (>60%), only 2 or 3 of them were used for each experiment. Concordance of viral dynamics among them, however, would give credence to the reproducibility of obtained results.

HBV/G infected ChiM mice by itself in corroboration with its monoinfection in human beings.²¹ The replication was very slow, however, and did not elevate serum HBV DNA to levels detectable by the method used (>10³ copies/ml). Coinfection with HBV/A enhanced the replication of HBV/G remarkably. HBV/G replicated vividly when coinfecting with HBV/C, as well. However, the time required for a 10-fold increase (log time) is 2-fold longer in mice initially infected with HBV/C versus HBV/A (3.3 versus 1.6 weeks). Combined, these results would indicate that HBV/G can thrive at the expense of other genotypes, and coinfection with HBV/A is much more advantageous for its enhanced replication than the other genotypes, including HBV/C. In support of this view, coinfection with HBV/A is frequent in individuals infected with HBV/G.^{16,34} Such a heavy dependence of HBV/G on HBV/A does not require recombination between them, because no recombination events occurred in ChiM mice coinfecting with them.

The initial replication of HBV/G was much slower than that of HBV/A, even in simultaneous coinfection. This was typically observed in three ChiM mice inoculated with a mouse passage of G-on-A superinfection containing HBV/G in the concentration a few logs higher than that of HBV/A. Despite such an enormous difference in introduced virions, the replication of HBV/A far

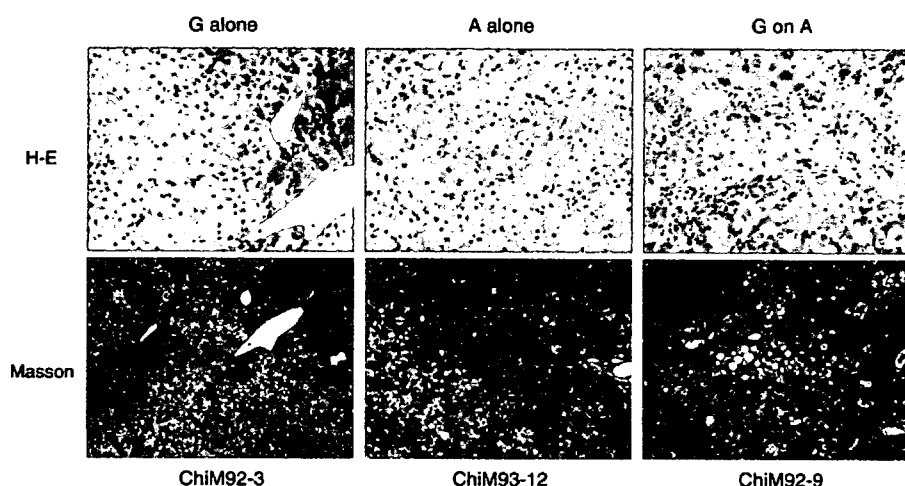


Fig. 9. Liver histology in a ChiM92-3 mouse monoinfected with genotype G, a ChiM93-12 mouse coinfecting with genotypes A and G (but persistently infected with genotype A alone), and a ChiM92-9 mouse superinfected with genotype G-on-A. Liver sections stained with hematoxylin-eosin or Masson's trichrome stain are shown.

exceeded that of HBV/G in the initial several weeks. Thereafter, HBV/G caught up with HBV/A, then took it over almost completely. Such a replacement was observed when HBV/G was superinfected on HBV/A, or vice versa.

The mechanism by which genotype G depends on genotype A for replication was pursued in cotransfection experiments in Huh7 cells. Cotransfection with the pcDNA_core clone carrying the core gene of genotype A2 increased the replication of the pUC19/G clone constructed with 1.24-fold the genome of genotype G in a dose-dependent manner (Fig. 8). Hence, *trans*-complementation with the core protein of genotype A would be required for genotype G to replicate actively. The possi-

bility remains, however, for other viral elements from coinfecting genotypes to enhance the replication of genotype G. Cotransfection of cultured cells with genotype G and others would help clarify how it depends on coinfecting the other genotypes.³⁵

Coinfection with HBV/G may be associated with pathological manifestations. ChiM mice coinfecting with HBV/A and HBV/G developed steatosis and fibrosis in the liver not observed in mice monoinfected with either of these genotypes. Very recently, Lacombe et al.²² reported more severe fibrosis in HIV-positive French patients who were infected with HBV/G than the others; they would most likely have been coinfecting with other genotypes in undetectable levels. On the basis of clinical and experimental pieces of evidence, it does seem that HBV/G has a strong disease-inducing capacity, which would be operable only when it is coinfecting with other genotypes. High levels of HBcrAg in mice with HBV/G (Figs. 3-6) under immunocompromised states would implicate accumulation of the product of the core gene in the fibrosis of patients coinfecting with it and HIV. Patients with HIV are infected with HBV at a frequency of 6%-9%, and liver-related deaths happen more often in coinfecting patients.^{36,37} Fibrosis proceeds faster in patients coinfecting with HIV and HBV, as in those with HCV.^{38,39} Therapeutic intervention to prevent fibrosis would be required in patients coinfecting with HIV and HBV, particularly in HBV/G patients.

In conclusion, the early viral dynamics of HBV/G have been characterized in ChiM mice monoinfected with HBV/G or coinfecting with other genotypes. The replication of HBV/G is very slow and depends heavily on coinfection with other genotypes. HBV/G rapidly takes over

Table 1. Steatosis and Fibrosis in Human Hepatocytes in the Liver of Chimeric Mice Monoinfected or Coinfecting with HBV/G

Inoculation	Mouse No.	Features	
		Steatosis (%) [*]	Fibrosis Stage
G alone	ChiM92-3	<5	F0
	ChiM 184-4	<5	F0
A alone	ChiM 93-12†	<5	F0
A-on-G	ChiM 93-4	50	F1
	ChiM 172-3	40	F1
G-on-A	ChiM 92-9	40	F2
	ChiM 124-11	50	F1
G-on-C	ChiM 91-21	80	F2
	ChiM 95-11	NA	NA
A plus G	ChiM 93-10	30	F0
Passage	ChiM 169-8	50	F1
	A plus G	ChiM 133-3	<5
	ChiM 133-6	30	F2

Abbreviation: NA, not available.

^{*}Percentage of human hepatocytes with steatosis. [†]Simultaneously inoculated with A plus G but became infected with genotype A only (Fig. 6B).

the other genotypes, though they are indispensable. Infection with HBV/G may induce steatosis and fibrosis in the liver—but again, only in the case of coinfection with other genotypes. However, it is still unclear whether or not such an increased pathogenicity of HBV/G is expressed exclusively in animals and patients with genetic or acquired immune deficiency.

References

- Lee WM. Hepatitis B virus infection. *N Engl J Med* 1997;337:1733-1745.
- Arauz-Ruiz P, Norder H, Robertson BH, Magnus LO. Genotype H: a new Amerindian genotype of hepatitis B virus revealed in Central America. *J Gen Virol* 2002;83:2059-2073.
- Norder H, Hammes B, Lofdahl S, Courouce AM, Magnus LO. Comparison of the amino acid sequences of nine different serotypes of hepatitis B surface antigen and genomic classification of the corresponding hepatitis B virus strains. *J Gen Virol* 1992;73:1201-1208.
- Okamoto H, Tsuda F, Sakugawa H, Sastrosoewignjo RI, Imai M, Miyakawa Y, et al. Typing hepatitis B virus by homology in nucleotide sequence: comparison of surface antigen subtypes. *J Gen Virol* 1988;69:2575-2583.
- Stuyver L, De Gendt S, Van Geyt C, Zoulim F, Fried M, Schinazi RF, et al. A new genotype of hepatitis B virus: complete genome and phylogenetic relatedness. *J Gen Virol* 2000;81:67-74.
- Naumann H, Schaefer S, Yoshida CF, Gaspar AM, Repp R, Gerlich WH. Identification of a new hepatitis B virus (HBV) genotype from Brazil that expresses HBV surface antigen subtype adw4. *J Gen Virol* 1993;74:1627-1632.
- Miyakawa Y, Mizokami M. Classifying hepatitis B virus genotypes. *Intervirology* 2003;46:329-338.
- Schaefer S. Hepatitis B virus: significance of genotypes. *J Viral Hepat* 2005;12:111-124.
- Chu CJ, Lok AS. Clinical significance of hepatitis B virus genotypes. *HEPATOLOGY* 2002;35:1274-1276.
- Liu CJ, Kao JH, Chen DS. Therapeutic implications of hepatitis B virus genotypes. *Liver Int* 2005;25:1097-1107.
- Sugauchi F, Kumada H, Sakugawa H, Komatsu M, Niitsuma H, Watanabe H, et al. Two subtypes of genotype B (Ba and Bj) of hepatitis B virus in Japan. *Clin Infect Dis* 2004;38:1222-1228.
- Tanaka Y, Orito E, Yuen MF, Mukaide M, Sugauchi F, Ito K, et al. Two subtypes (subgenotypes) of hepatitis B virus genotype C: A novel subtyping assay based on restriction fragment length polymorphism. *Hepatol Res* 2005;33:216-224.
- Sugauchi F, Kumada H, Acharya SA, Shrestha SM, Gamutan MT, Khan M, et al. Epidemiological and sequence differences between two subtypes (Ae and Aa) of hepatitis B virus genotype A. *J Gen Virol* 2004;85:811-820.
- Akuta N, Suzuki F, Kobayashi M, Tsubota A, Suzuki Y, Hosaka T, et al. The influence of hepatitis B virus genotype on the development of lamivudine resistance during long-term treatment. *J Hepatol* 2003;38:315-321.
- Tanaka Y, Hasegawa I, Kato T, Orito E, Hirashima N, Acharya SK, et al. A case-control study for differences among hepatitis B virus infections of genotypes A (subtypes Aa and Ae) and D. *HEPATOLOGY* 2004;40:747-755.
- Kato H, Orito E, Gish RG, Sugauchi F, Suzuki S, Ueda R, et al. Characteristics of hepatitis B virus isolates of genotype G and their phylogenetic differences from the other six genotypes (A through F). *J Virol* 2002;76:6131-6137.
- Kato H, Orito E, Gish RG, Bzowej N, Newsom M, Sugauchi F, et al. Hepatitis B e antigen in sera from individuals infected with hepatitis B virus of genotype G. *HEPATOLOGY* 2002;35:922-929.
- Perez-Olmeda M, Nunez M, Garcia-Samaniego J, Rios P, Gonzalez-Lahoz J, Soriano V. Distribution of hepatitis B virus genotypes in HIV-infected patients with chronic hepatitis B: therapeutic implications. *AIDS Res Hum Retroviruses* 2003;19:657-659.
- Suwannakarn K, Tangkijvanich P, Theamboonlers A, Abe K, Poorawan Y. A novel recombinant of hepatitis B virus genotypes G and C isolated from a Thai patient with hepatocellular carcinoma. *J Gen Virol* 2005;86:3027-3030.
- Sanchez LV, Tanaka Y, Maldonado M, Mizokami M, Panduro A. Difference of hepatitis B virus genotype distribution in two groups of Mexican patients with different risk factors. High prevalence of genotype H and G. *Intervirology* 2007;50:9-15.
- Chudy M, Schmidt M, Czudai V, Scheiblaue H, Nick S, Mosebach M, et al. Hepatitis B virus genotype G mono-infection and its transmission by blood components. *HEPATOLOGY* 2006;44:99-107.
- Lacombe K, Massari V, Girard PM, Serfaty L, Gozlan J, Pialoux G, et al. Major role of hepatitis B genotypes in liver fibrosis during coinfection with HIV. *AIDS* 2006;20:419-427.
- Heckel JL, Sandgren EP, Degen JL, Palmiter RD, Brinster RL. Neonatal bleeding in transgenic mice expressing urokinase-type plasminogen activator. *Cell* 1990;62:447-456.
- Rhim JA, Sandgren EP, Degen JL, Palmiter RD, Brinster RL. Replacement of diseased mouse liver by hepatic cell transplantation. *Science* 1994;263:1149-1152.
- Tateno C, Yoshizane Y, Saito N, Kataoka M, Utoh R, Yamasaki C, et al. Near completely humanized liver in mice shows human-type metabolic responses to drugs. *Am J Pathol* 2004;165:901-912.
- Mercer DF, Schiller DE, Elliott JF, Douglas DN, Hao C, Rinfret A, et al. Hepatitis C virus replication in mice with chimeric human livers. *Nat Med* 2001;7:927-933.
- Tsuge M, Hirano N, Takaishi H, Noguchi C, Oga H, Imamura M, et al. Infection of human hepatocyte chimeric mouse with genetically engineered hepatitis B virus. *HEPATOLOGY* 2005;42:1046-1054.
- Sugiyama M, Tanaka Y, Kato T, Orito E, Ito K, Acharya SK, et al. Influence of hepatitis B virus genotypes on the intra- and extracellular expression of viral DNA and antigens. *HEPATOLOGY* 2006;44:915-924.
- Kimura T, Ohno N, Terada N, Rokuhara A, Matsumoto A, Yagi S, et al. Hepatitis B virus DNA-negative Dane particles lack core protein but contain a 22-kDa precore protein without C-terminal arginine-rich domain. *J Biol Chem* 2005;280:21713-21719.
- Shinkai N, Tanaka Y, Orito E, Ito K, Ohno T, Hirashima N, et al. Measurement of hepatitis B virus core-related antigen as predicting factor for relapse after cessation of lamivudine therapy for chronic hepatitis B virus infection. *Hepatol Res* 2006;36:272-276.
- Abe A, Inoue K, Tanaka T, Kato J, Kajiyama N, Kawaguchi R, et al. Quantitation of hepatitis B virus genomic DNA by real-time detection PCR. *J Clin Microbiol* 1999;37:2899-2903.
- Mason AL, Xu L, Guo L, Kuhns M, Perrillo RP. Molecular basis for persistent hepatitis B virus infection in the liver after clearance of serum hepatitis B surface antigen. *HEPATOLOGY* 1998;27:1736-1742.
- Kato H, Orito E, Sugauchi F, Ueda R, Gish RG, Usuda S, et al. Determination of hepatitis B virus genotype G by polymerase chain reaction with hemi-nested primers. *J Virol Methods* 2001;98:153-159.
- Osiowy C, Giles E. Evaluation of the INNO-LiPA HBV genotyping assay for determination of hepatitis B virus genotype. *J Clin Microbiol* 2003;41:5473-5477.
- Kremsdorff D, Garreau F, Capel F, Petit MA, Brechot C. In vivo selection of a hepatitis B virus mutant with abnormal viral protein expression. *J Gen Virol* 1996;77:929-939.
- Konopnicki D, Mocroft A, de Wit S, Antunes F, Ledergerber B, Katlama C, et al. Hepatitis B and HIV: prevalence, AIDS progression, response to highly active antiretroviral therapy and increased mortality in the EuroSIDA cohort. *AIDS* 2005;19:593-601.
- Thio CL, Seaberg EC, Skolasky R Jr, Phair J, Visscher B, Munoz A, et al. HIV-1, hepatitis B virus, and risk of liver-related mortality in the Multicenter Cohort Study (MACS). *Lancet* 2002;360:1921-1926.
- Benhamou Y, Bochet M, Di Martino V, Charlotte F, Azria F, Coutellier A, et al. Liver fibrosis progression in human immunodeficiency virus and hepatitis C virus coinfecting patients. The Multivirc Group. *HEPATOLOGY* 1999;30:1054-1058.
- Colin JF, Cazals-Hatem D, Lioriot MA, Martinot-Peignoux M, Pham BN, Auperin A, et al. Influence of human immunodeficiency virus infection on chronic hepatitis B in homosexual men. *HEPATOLOGY* 1999;29:1306-1310.

Diversity of preferred nucleotide sequences around the translation initiation codon in eukaryote genomes

So Nakagawa¹, Yoshihito Niimura², Takashi Gojobori^{3,4}, Hiroshi Tanaka^{1,2,*} and Kin-ichiro Miura⁵

¹Department of Systems Biology, School of Biomedical Science, ²Department of Bioinformatics, Medical Research Institute, Tokyo Medical and Dental University, Yushima, Tokyo, ³Center for Information Biology and DNA Data Bank of Japan, National Institute of Genetics, Mishima, Shizuoka, ⁴Biological Information Research Center, National Institute of Advanced Industrial Science and Technology, Aomi, Tokyo and ⁵Department of Medical Genome Science, Graduate School of Frontier Science, University of Tokyo, Kashiwa, Chiba, Japan

Received August 21, 2007; Revised November 2, 2007; Accepted November 27, 2007

ABSTRACT

Understanding regulatory mechanisms of protein synthesis in eukaryotes is essential for the accurate annotation of genome sequences. Kozak reported that the nucleotide sequence **GCCGCC (A/G)CCAUGG** (**AUG** is the initiation codon) was frequently observed in vertebrate genes and that this 'consensus' sequence enhanced translation initiation. However, later studies using invertebrate, fungal and plant genes reported different 'consensus' sequences. In this study, we conducted extensive comparative analyses of nucleotide sequences around the initiation codon by using genomic data from 47 eukaryote species including animals, fungi, plants and protists. The analyses revealed that preferred nucleotide sequences are quite diverse among different species, but differences between patterns of nucleotide bias roughly reflect the evolutionary relationships of the species. We also found strong biases of **A/G** at position **-3**, **A/C** at position **-2** and **C** at position **+5** that were commonly observed in all species examined. Genes with higher expression levels showed stronger signals, suggesting that these nucleotides are responsible for the regulation of translation initiation. The diversity of preferred nucleotide sequences around the initiation codon might be explained by differences in relative contributions from two distinct patterns, **GCCGCCAUG** and **AAAAAAUG**, which implies the presence of multiple molecular mechanisms for controlling translation initiation.

INTRODUCTION

The control of translation initiation is one of the most fundamental processes in the regulation of gene expression. In 1978, Kozak (1,2) proposed the scanning model for translation initiation in eukaryotes. According to this model, the 40S ribosomal subunit with several initiation factors binds the 7-methyl guanosine cap at the 5' end of an mRNA and moves along the mRNA until it encounters an AUG codon. It was also proposed that when the AUG codon is in the context of **GCCGCC(A/G)CCAUGG** (**A/G** represents A or G and **AUG** represents the translation initiation codon), which is called the 'Kozak consensus sequence', the efficiency of translation initiation is enhanced. However, the detailed molecular mechanism of translation initiation in eukaryotes is still unclear. Moreover, although the sequence is described as a 'consensus' sequence, the extent of conservation is quite low. It was reported that only 0.2% of vertebrate genes contain precisely the sequence **GCCGCC(A/G)CCAUGG** (3). We therefore avoid using the word 'consensus' in this context, and instead refer to the sequence as 'preferred' sequence.

Kozak compiled 211 genes (4) and 699 genes (5) primarily from vertebrates and obtained the above sequence. This sequence was initially thought to be essential for all eukaryotes (4). Later, however, it was revealed that a preferred nucleotide sequence around the initiation codon varies considerably among different species. The preferred sequences are **GCGGC(A/C)(A/G)(A/C)CAUGGCG** for Monocots (1127 genes), **AAAAAAA(A/C)AAUGGCU** for Dicots (derived from 3643 genes) (6), **ACAACCAAAAUGGC** for *Drosophila melanogaster* (192 genes), **UAAAT(A/C)AACCAUG(A/G)C** for other invertebrates (155 genes), and

*To whom correspondence should be addressed. Tel: +81 3 5803 5839; Fax: +81 3 5803 0247; Email: htanaka@bioinfo.tmd.ac.jp

© 2007 The Author(s)

This is an Open Access article distributed under the terms of the Creative Commons Attribution Non-Commercial License (<http://creativecommons.org/licenses/by-nc/2.0/uk/>) which permits unrestricted non-commercial use, distribution, and reproduction in any medium, provided the original work is properly cited.

AAAAAAAAAUGTC for *Saccharomyces cerevisiae* (461 genes) (3). Kozak also reported that replacement of A/G at position -3 (three bases before the initiation codon) and G at position $+4$ (one base after the initiation codon) strongly impaired translation initiation in mammals (7,8). However, in *S. cerevisiae* nucleotide substitutions at position -3 did not substantially affect the rate of translation initiation (9,10), although there is a nucleotide bias towards A at this position (3). It therefore appears that the molecular mechanisms for recognizing the initiation codon vary among species.

There have been two limitations to previous studies aimed at identifying preferred sequences around the initiation codon. First, the number of species and genes examined was limited. In this study, we used whole-genome expression data and gene sequences from diverse eukaryote species. The second issue has been that the GC contents in genomes are known to differ from species to species. The preference for A before the initiation codon in Dicots and *S. cerevisiae* can be partially explained by the AT-richness of their genomes. To compare nucleotide sequences responsible for translation initiation among various species, differences in the usage of nucleotides in each genome must be considered. We previously invented a method of graphically representing nucleotide appearance biases at each position in a gene on the basis of the deviation from the expected values that are calculated for a given genome sequence (11,12). Application of this method to bacterial genomic data led to the successful identification of the Shine-Dalgarno (SD) sequence, a well-characterized signal for translation initiation in prokaryotes (11). We have also reported that the nucleotides appearing at the second codon (the codon next to the initiation codon) are highly biased in eukaryote genes and that a preferred second codon is characteristic of each species (e.g. GCG for mammals and plants) (12).

To obtain additional insight into the molecular mechanisms of translation initiation in eukaryotes, we extensively examined the nucleotide sequences around the initiation codon by using the method introduced above. We conducted comparative analyses of the biases in nucleotides located in positions proximal to the initiation codon among 47 eukaryote species including animals, fungi, plants and protists. We thereby were able to identify both universal and species-specific features, and these features possibly reflect the evolution of the mechanism of translation initiation.

MATERIALS AND METHODS

Data

We used cDNA or genome sequence data from 47 eukaryote species including 22 metazoans, eight plants, nine fungi and eight protists. Species names and the database used are shown in Table 1. We used only protein-coding genes that start from the AUG codon and end with a stop codon. As for human genes, we used genes in categories I–IV provided by the H-Invitational Database (13). When information about alternative splicing variants was available, only one representative sequence with the

longest coding sequence (CDS) was used. Otherwise, all of the protein-coding genes were used [for UniGene database (14)]. The amount of expressed mRNAs in humans and *S. cerevisiae*, obtained by serial analysis of gene expression (SAGE), were downloaded from H-ANGEL (<http://jbirc.jbic.or.jp/hinv/h-angel/>) (15) and Holstege's web site (<http://www.wi.mit.edu/young/expression.html>) (16), respectively.

Evaluation of nucleotide frequency bias

To examine biases in nucleotide appearance around the initiation codon, all genes from each species were aligned at the initiation codons without any alignment gaps. The number of each nucleotide [A, U (T), G and C] was counted at each position in the alignment. The observed numbers of nucleotides were compared with the expected numbers using the likelihood-ratio statistic or the G -statistic, which is used for a test for goodness-of-fit (17). The expectations were calculated for each species in four separate categories, namely, the 5' untranslated regions (UTRs) and the first, second and third positions in a codon in CDSs, because nucleotide frequencies are different among these categories. The G -value at position i was calculated by the formula:

$$G^{(i)} = 2 \sum_n O_n^{(i)} \left(\frac{O_n^{(i)}}{E_n^{(i)}} \right) \quad 1$$

where $O_n^{(i)}$ is the observed number of nucleotide n (A, U, G and C) at position i , and $E_n^{(i)}$ is the expected number of nucleotide n in the category to which position i belongs (5' UTRs or the first, second or third positions in a codon). As regards the genomic data [RefSeq, MIPS and GeneDB (14,18,19)], 100 base-pair (bp) regions upstream from the initiation codon were regarded as the 5' UTRs and data from these regions were used for the computation of the expectations. It is known that the distribution of the G -statistic is approximated by the χ^2 -distribution with $f-1$ degrees of freedom when the sample size is large, where f is the number of different classes ($f=4$). Each term in Formula 1 represents the contribution of each nucleotide to the bias. When $O_n^{(i)}$ is larger and smaller than $E_n^{(i)}$, the values of $2O_n^{(i)} \ln(O_n^{(i)}/E_n^{(i)})$ become positive and negative, respectively. For this reason, we regarded each term in Formula 1 as a measure of the bias for each nucleotide at a given position. G -values are proportional to the number of genes (N) when the fractions of observed and expected numbers of nucleotides are the same. To compare nucleotide biases among different species with different numbers of genes, we defined a value that is not affected by the number of genes, $g_n^{(i)} = 2o_n^{(i)} \ln(o_n^{(i)}/e_n^{(i)})$, where $o_n^{(i)}$ and $e_n^{(i)}$ are the fractions of the observed and expected numbers of nucleotide n at position i . When $o_n^{(i)}$ is zero, $g_n^{(i)}$ is defined to be zero. The G -value divided by N is equal to the sum of $g_n^{(i)}(G^{(i)}/N = \sum_n g_n^{(i)})$.

Cluster analysis of the patterns in nucleotide biases

We quantified similarities between the patterns in nucleotide bias around initiation codons by using the Pearson's correlation coefficient. The correlation coefficient r_{XY}

Table 1. The 47 eukaryote species used for analysis

Species	Common name	Database ^a
Animals, Vertebrates		
<i>Homo sapiens</i> ^b	Human	H-Invitational Database 3.0 (13)
<i>Pan troglodytes</i> ^b	Chimpanzee	Ensembl (CHIMPIA) (35)
<i>Macaca fascicularis</i> ^b	Crab-eating macaque	UniGene (14)
<i>Macaca mulatta</i> ^b	Rhesus monkey	Ensembl (MMUL_0_1)
<i>Mus musculus</i> ^b	Mouse	FANTOM3 (36)
<i>Rattus norvegicus</i> ^b	Rat	Mammalian Gene Collection (37)
<i>Oryctolagus cuniculus</i> ^b	Rabbit	UniGene
<i>Canis familiaris</i> ^b	Dog	Ensembl (BROADD1)
<i>Bos taurus</i> ^b	Cattle	Mammalian Gene Collection
<i>Sus scrofa</i> ^b	Pig	UniGene
<i>Gallus gallus</i> ^b	Chicken	Ensemble (WASHUC1)
<i>Xenopus laevis</i> ^b	African clawed frog	Xenopus Gene Collection (38)
<i>Xenopus tropicalis</i> ^b	Western clawed frog	Xenopus Gene Collection
<i>Danio rerio</i> ^b	Zebrafish	Zebrafish Gene Collection (39)
Animals, Invertebrates		
<i>Ciona intestinalis</i> ^b	Sea squirt	UniGene
<i>Drosophila melanogaster</i> ^b	Fruit fly	Ensemble (BDGP4)
<i>Anopheles gambiae</i> ^b	African malaria mosquito	Ensemble (AgamP3)
<i>Apis mellifera</i> ^b	Honeybee	Ensemble (AMEL2.0)
<i>Bombyx mori</i> ^b	Domestic silkworm	UniGene
<i>Tribolium castaneum</i>	Red flour beetle	RefSeq (14)
<i>Caenorhabditis elegans</i> ^b		Ensemble (CEL150)
<i>Schistosoma japonicum</i> ^b		UniGene
Plants, Monocots		
<i>Oryza sativa</i> ^b	Rice	KOME (released on 24 December 2004) (40)
<i>Hordeum vulgare</i> ^b	Barley	UniGene
<i>Triticum aestivum</i> ^b	Bread wheat	UniGene
<i>Zea mays</i> ^b	Indian corn	UniGene
Plants, Dicots		
<i>Arabidopsis thaliana</i> ^b	Thale cress	TAIR (released on 28 February 2004) (41)
<i>Glycine max</i> ^b	Soybean	UniGene
<i>Lycopersicon esculentum</i> ^b	Tomato	UniGene
<i>Solanum tuberosum</i> ^b	Potato	UniGene
Fungi		
<i>Saccharomyces cerevisiae</i>	Budding yeast	MIPS (18)
<i>Debaryomyces hansenii</i>		RefSeq
<i>Eremothecium gossypii</i>		RefSeq
<i>Kluyveromyces lactis</i>		RefSeq
<i>Yarrowia lipolytica</i>		RefSeq
<i>Candida glabrata</i>		RefSeq
<i>Schizosaccharomyces pombe</i>	Fission yeast	GeneDB (Version 2.1) (19)
<i>Aspergillus fumigatus</i>		RefSeq
<i>Cryptococcus neoformans</i>		RefSeq
Protists		
<i>Theileria parva</i>		RefSeq
<i>Theileria annulata</i>		RefSeq
<i>Cryptosporidium parvum</i>		RefSeq
<i>Plasmodium falciparum</i>		GeneDB (released on 26 January 2006)
<i>Leishmania major</i>		RefSeq
<i>Trypanosoma brucei</i>		RefSeq
<i>Dictyostelium discoideum</i>	Slime mold	dictyBase (released on 3, May, 2006) (42)
<i>Cyanidioschyzon merolae</i>		<i>Cyanidioschyzon merolae</i> Genome Project (43)

^aThese data were downloaded from the following websites. H-Invitational Database 3.0, <http://www.jbirc.jbic.or.jp/hinv/ahg-db/>; Ensembl, <http://www.ensembl.org/>; UniGene, <http://www.ncbi.nlm.nih.gov/sites/entrez?db=unigene>; FANTOM3, <http://fantom.gsc.riken.go.jp/>; Mammalian Gene Collection, <http://mgc.nci.nih.gov/>; Xenopus Gene Collection, <http://xgc.nci.nih.gov/>; Zebrafish Gene Collection, <http://zgc.nci.nih.gov/>; RefSeq, <http://www.ncbi.nlm.nih.gov/RefSeq/>; KOME, <http://cdna01.dna.affrc.go.jp/cDNA/>; TAIR, <http://www.arabidopsis.org/>; MIPS, <http://mips.gsf.de/>; GeneDB, <http://www.sanger.ac.uk/>; dictyBase, <http://dictybase.org/>; *Cyanidioschyzon merolae* Genome Project, <http://merolae.biol.s.u-tokyo.ac.jp/>.

^bSpecies for which cDNA data were used.

between species X and Y was calculated from the g_n values from positions -9 to -1 in the 5' UTRs, and from positions $+4$ to $+6$ in the CDSs (the second codon) as follows:

$$r_{XY} = \frac{\sum_i \sum_n (g_{Xn}^{(i)} - \bar{g}_X)(g_{Yn}^{(i)} - \bar{g}_Y)}{\sqrt{\sum_i \sum_n (g_{Xn}^{(i)} - \bar{g}_X)^2} \sqrt{\sum_i \sum_n (g_{Yn}^{(i)} - \bar{g}_Y)^2}}$$

where $g_{Xn}^{(i)}$ and $g_{Yn}^{(i)}$ represent g_n values of nucleotide n (A, U, G or C) at position i in species X and Y, respectively, and \bar{g}_X and \bar{g}_Y represent the average of g_n values among all positions (from -9 to -1 and from $+4$ to $+6$) and nucleotides in species X and Y, respectively. We calculated r -values for all combinations among the 47 species examined and defined the similarity score D as $1 - r$. Using the similarity scores, the cluster analysis was conducted by the group average method (Figure 3), the centroid method and the Ward method (Figure S1). Note that D is free from the absolute values of g_n . When the number of genes used is small, the g_n values tend to become large, apparently because highly expressed genes are more likely to be contained in a small gene set than are genes expressed at low levels. Therefore, although the g_n values are affected by the number of genes, the values of D are expected to be robust against the difference in the number of genes for each species.

Evaluation of hexanucleotide biases

We evaluated the deviation of the observed number of a particular combination of nucleotides from the expected number by using the Z -value (Tables 3 and 4 and Table S1 that is available as Supplementary Data). The Z -value was calculated by $Z = (O - E)/[E(1 - E/N)]^{1/2}$, where N is the number of genes, and O and E are observed and expected numbers of a particular combination of nucleotides, respectively. The expected number was calculated by assuming that a nucleotide at each position appears independently (see the legend of Table 3). We calculated the Z -values for all possible combinations of six nucleotides ($4^6 = 4096$ combinations) in the region upstream of the initiation codon (from positions -6 to -1) and ranked them according to the Z -value. To see whether a sequence that is a mixture of GCCGCCAUG and AAAAAAUG is suppressed in genes or not, we examined hexanucleotide sequences generated by combining three nucleotides from GCCGCCAUG and three nucleotides from AAAAAAUG (e.g. GAAACCAUG or ACAGACAUG). There are 20 (6C_3) such combinations (Table 4). From them, we excluded AAAGCCAUG and GCCAAAUG, because GCCAUG and AAAAUG were observed much more frequently than the expectations (Table S1). We regarded the remaining 18 sequences as 'mixed sequences'. We conducted the Wilcoxon rank sum test to see whether the ranks of the 18 mixed sequences are significantly low among the 4096 sequences or not.

RESULTS

We evaluated biases in nucleotide appearance at each position around the initiation codon by using the

G -statistic (see Materials and Methods section). Figure 1 shows the results obtained for 10 012 human genes. As shown in the upper diagram of this figure, the fractions of nucleotides A, T, G and C vary considerably in a position-dependent manner. The largest deviation of nucleotide frequencies from the expected values was observed at position -3 , which is indicated by the largest G -value at this position (middle diagram). At this position, the values of g_n are positive for A and G (lower diagram), which indicates that A and G appear more frequently than the expectations (see Materials and Methods section). In fact, the fractions of A and G at position -3 are 39.3% and 34.6%, respectively, which are much larger than those in the entire 5' UTRs of 10 012 human genes (23.7% and 26.6%, respectively). The results depicted in the lower diagram suggest that the preferred sequence in humans is GCCGCC(A/G)(C/A)CAUGGCG, which is nearly the same as the sequence reported by Kozak (4,5). Note that the bias of GCG at the second codon is also quite strong, as we previously reported (12).

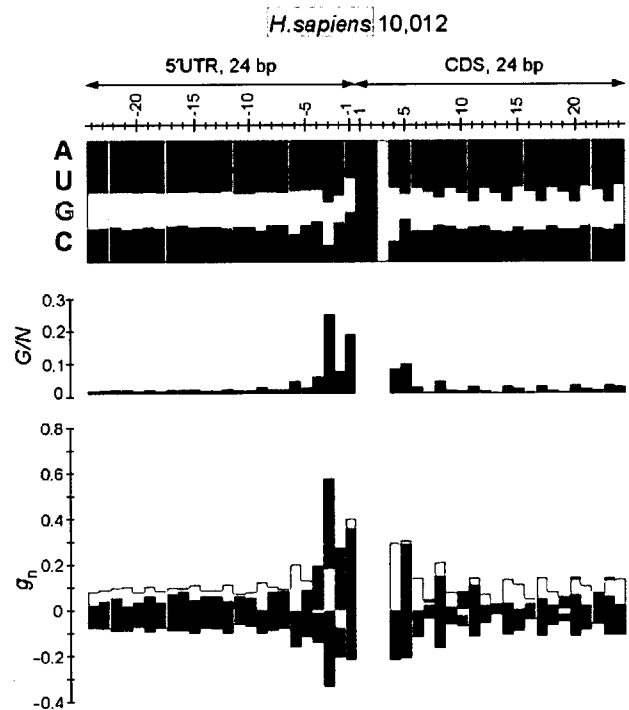


Figure 1. Biases in nucleotide appearance for 10012 human genes. Top, fractions of nucleotides appearing at each position in 24 base-pair (bp) regions in 5' UTRs and CDSs. A, U, G and C are shown in green, magenta, yellow and blue, respectively. Middle, G -values divided by the number of genes ($N = 10012$), showing the deviation from the expected values. Bottom, the values of g_n for $n = A, U, G$ or C at each position. The color scheme is the same as that used in the diagram at the top. Colored bars above and below the horizontal line indicate positive and negative g_n values, respectively, and these bars were drawn without overlapping. In the middle and bottom diagrams, the values for the initiation codon (AUG) are omitted. Note that the biases shown in this figure are statistically highly significant ($P < 10^{-10}$ from positions -9 to $+6$) because of the large sample size.

We applied the method described here to investigate the genes of 47 eukaryote species for which full-length cDNA or whole-genome data are available (Figure 2). These species included a wide variety of eukaryotes such as the soil-dwelling social amoeba *Dictyostelium discoideum* and the unicellular red alga *Cyanidioschyzon merolae* (Table 1). We used cDNA data when they are available, because gene annotation based on expression data is expected to be more accurate than that predicted from genome sequences. For most of the animal and plant species examined, cDNA data were used (Table 1). In Figure 2, only the region from positions -9 to +6 is shown, in which the G -values are relatively large. This figure reveals that the preferred nucleotide sequences around the initiation codon, as well as the extent of deviation from the expectations, vary among species. However, several features were commonly observed among species. For example, A is preferred at position -3 in all species examined. To compare the patterns of bias in nucleotide frequencies among different species, we quantified the similarity of g_n values in the region from positions -9 to +6 between two species (see Materials and Methods section). By using a similarity score (D), we conducted the cluster analysis (Figure 3). The results showed that vertebrates, Monocots and Dicots each formed a cluster, thus indicating that the patterns of nucleotide bias are similar within each of these groups of organisms. Although the cluster dendrogram changed depending on the method of cluster analysis used, the clustering of vertebrates, Monocots and Dicots was robust (Figure S1). Fungi, invertebrates (containing urochordates, arthropods, nematodes and platyhelminthes) and each taxonomic group of protists also tended to form a cluster. These observations suggest that the patterns of nucleotide bias around the initiation codon roughly reflect the evolutionary relationships of eukaryote species.

Figure 4 shows the preferred sequences for each taxonomic group of eukaryotes. These sequences were obtained by taking the average of the patterns of nucleotide bias for all species belonging to each group. The sequences obtained here are similar to those previously reported. For example, for Monocots the preferred sequence obtained is G(A/C)(G/A)GC(A/C/G)(G/A)(C/A)(G/C)AUGGCG, which is similar to that reported in Joshi *et al.* (6) (see Introduction section). The following biases in nucleotide appearance were commonly observed among all taxonomic groups examined: -6G (G at position -6), -3A/G, -2A/C and +5C. Of these biases, -3A is the most remarkable. Moreover, a general tendency toward the under-representation of U around the initiation codon was observed. The biases in protists are relatively weak, reflecting highly variable patterns of nucleotide bias in these species (Figure 2).

Figure 4 also suggests that A-rich biases are present in the region from positions -4 to -1 in almost all species examined. These biases are clearly observed even in species with very low GC content. For example, the fraction of A in the *D. discoideum* genome is 38.8% (the GC content is 22.4%) (20), while the fraction of A at position -3 is as high as 85.9% (Figure S2). Moreover, we identified signals that had not been reported to date.

Monocots showed a signal of GC(C/G)GC(C/G)AUG as mentioned above, but a similar pattern was also observed in Dicots. Furthermore, a relatively weak but clear signal of GCCGCAUG was detected from invertebrates and fungi, which is similar to the sequence for vertebrates. It therefore appears that the preferred sequences in eukaryotes can be regarded as a summation of the repetition of GCC and that of A (see Discussion section).

To determine whether the biases described above are responsible for the efficiency of translation initiation, we examined the correlation of the strengths of the biases with gene expression levels. It is reasonable to assume that the translation rates would be high for highly expressed genes and low for genes expressed at low levels. In fact, we have conducted a genome-wide microarray analysis and shown that the efficiency of translation initiation is correlated with the expression level of mRNAs for *S. cerevisiae* genes (Akiyama *et al.*, unpublished data). Therefore, the signals for efficient translation initiation are assumed to be more conspicuous for highly expressed genes. Figure 5A indicates the results for 1000 genes with high expression levels and those for 1000 genes with low expression levels in humans and *S. cerevisiae*. Table 2 gives the fractions of -3A and +5G for genes expressed at high and low levels, and those for the entire set of genes in the two species. These results clearly show that the biases identified by using an entire set of genes became stronger when highly expressed genes were used. These results suggest that some of the preferred sequences identified in this study are responsible for the efficiency of translation initiation.

We also identified a clear pattern of three-base periodicity from several vertebrate and Monocot species (Figure 5B). Interestingly, a similar pattern of a GCC or GCG repeat was observed in both 5' UTRs and CDS regions, and the biases were more prominent in the regions near the initiation codon. One might suspect that the three-base periodicity in the 5' UTRs is due to an artifact, i.e. that the CDSs are wrongly annotated as UTRs because of an inaccurate assignment of initiation codons. To determine whether this observation is due to an artifact or not, we conducted the same analysis using only the genes containing an in-frame stop codon in the 5' UTR, but which do not contain any in-frame AUG codons between the initiation codon and the closest upstream stop codon from the initiation codon (Figure S3A). For such genes, there are no possibilities that wrongly annotated 5' UTRs used in the analyses contain CDSs. The results clearly show that the periodic pattern described above is also observed in such genes, suggesting that this pattern is not an artifact, but rather a signal for the initiation of translation (Figure S3B).

As shown in Figure 5C, U- and A-rich biases were commonly observed around positions -40 and -15, respectively, in amphibians, fishes and insects. In *D. melanogaster*, for example, the average fraction of U in the region from positions -45 to -35 is 28.2%, and the average fraction of A from positions -20 to -10 is 35.6%. These values are considerably larger than the averages in the entire 5' UTRs (21.4% for U and 31.5% for A). However, these biases are not clear in the other species.

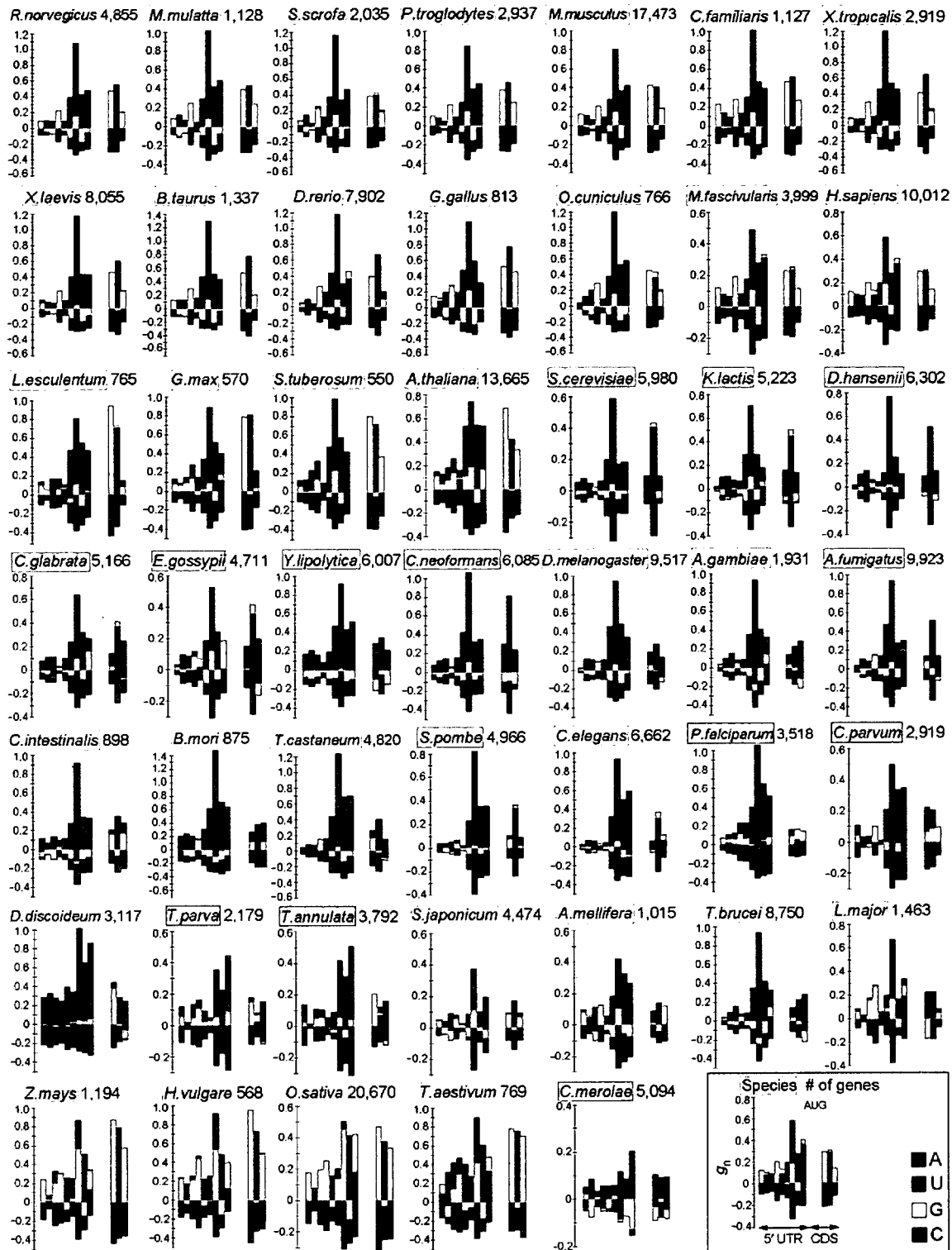


Figure 2. Nucleotide biases around the initiation codon in 47 eukaryote species. Each diagram shows g_n values from positions -9 to $+6$ in each species. The name of the species and the number of genes used are also given. The color scheme is the same as that used in Figure 1. The initiation codon (AUG) is not shown.

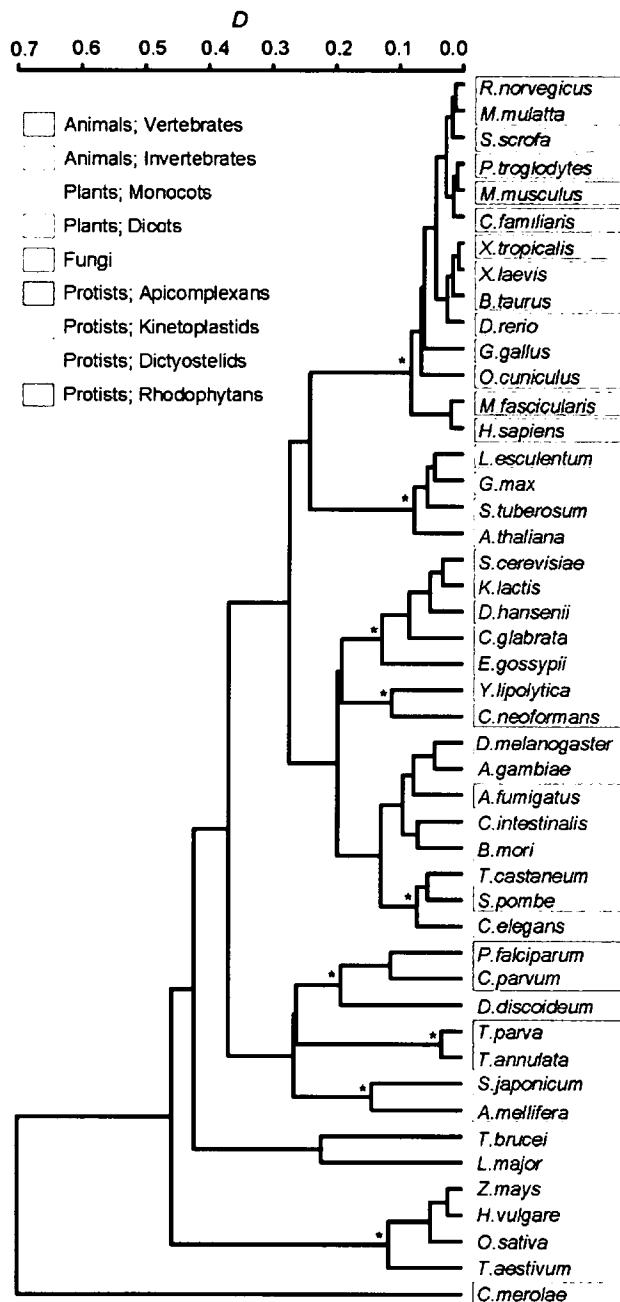


Figure 3. Cluster dendrogram showing similarities between patterns of nucleotide bias. The distance D was calculated from the g_n values from positions -9 to $+6$ (see Materials and Methods section). The group average method was used for the construction of the cluster dendrogram. An asterisk (*) indicates a cluster that is conserved among the dendrograms constructed by three different clustering methods (Figure S1).

DISCUSSION

In this study, we revealed that the signals $-3A/G$, $-2A/C$ and $+5C$ are common among various eukaryote species (Figure 4). Several studies have shown that $-3A/G$ plays

the most crucial role in enhancing translation initiation (7,21–24). In accord with these studies, our results indicated that the signal of $-3A/G$ is the most remarkable in almost all eukaryote species examined, and this signal is even stronger in highly expressed genes (Figure 5A). Recently, Pisarev *et al.* (24) demonstrated that $-3G$ in an mRNA interacts with a eukaryotic initiation factor eIF2 α by using a rabbit cell system, although the amino acids in an eIF2 α that are involved in the interaction are still unknown. Based on this observation, a purine base at position -3 was proposed to interact with an eIF2 α as a key element in translation initiation. Since the amino acid sequences of this protein are highly conserved among various eukaryotes, the interaction between the nucleotide at position -3 and an eIF2 α may be a common mechanism for translation initiation.

Although Kozak (25) reported that recognition of the initiation codon was not augmented by a nucleotide at position $+5$, other researchers suggested that $+5A/C$ in mammals or $+5C$ in plants affect the efficiency of translation initiation (26,27). Our analyses using various eukaryotes (Figures 2 and 4) and highly expressed genes (Figure 5A and Table 2) suggest the importance of $+5C$ for translation initiation. However, since the nucleotide at position $+5$ determines the chemical properties of the second amino acid, it is also possible that this nucleotide is under the functional constraint of the amino acid sequences (12). As regards the $-2A/C$ signal, to our knowledge, there has been no experimental data suggestive of its role in the initiation of translation.

Although $+4G$ has been described as important for translation initiation in vertebrates and plants, the effect of $+4G$ is relatively minor compared with that of $-3A/G$ (21,24,25). Our study showed that the nucleotide appearing at position $+4$ is highly biased, but a preferred nucleotide is not common among all eukaryotes. In vertebrates and plants, G is preferred at this position, whereas in invertebrates, fungi and protists, T is generally preferred. The biases of $+4G$ in humans and $+4U$ in *S. cerevisiae* were even more conspicuous when highly expressed genes were examined, suggesting the possibility that position $+4$ is involved in enhancing translation initiation; however, the nucleotide at the position required for effective translation initiation appear to be diversified among eukaryotes.

In the original scanning model, it was postulated that translation is initiated from the first AUG codon in an mRNA (8). However, it has since been revealed that the actual mechanism of translation initiation is much more complicated than previously thought. It is known that AUG trinucleotides referred to as upstream AUGs (uAUGs) are frequently observed in 5' UTRs, and that short open reading frames designated as upstream ORFs (uORFs) are often also present (28). It has been reported that $\sim 55\%$ and $\sim 25\%$ of mammalian genes have one or more uAUGs and uORFs, respectively (29). These uAUGs and uORFs are apparently involved in the down-regulation of translation (30). Moreover, even if the first AUG codon is located within a context of a 'Kozak consensus sequence', translation is not necessarily initiated from it (31,32). Dresios *et al.* (33) suggested that a short element in a

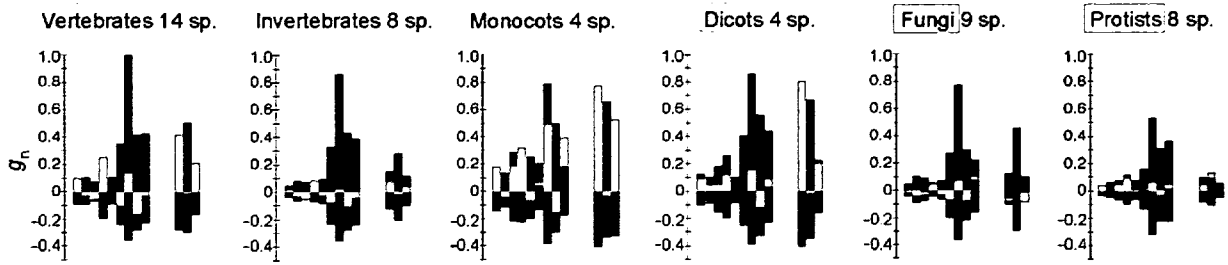


Figure 4. Nucleotide biases around the initiation codon for each taxonomic group of eukaryotes. The g_n value at each position for each nucleotide was calculated from the average of o_n values and that of e_n values among all species belonging to a given taxonomic group. The number of species used is shown for each group. sp., species.

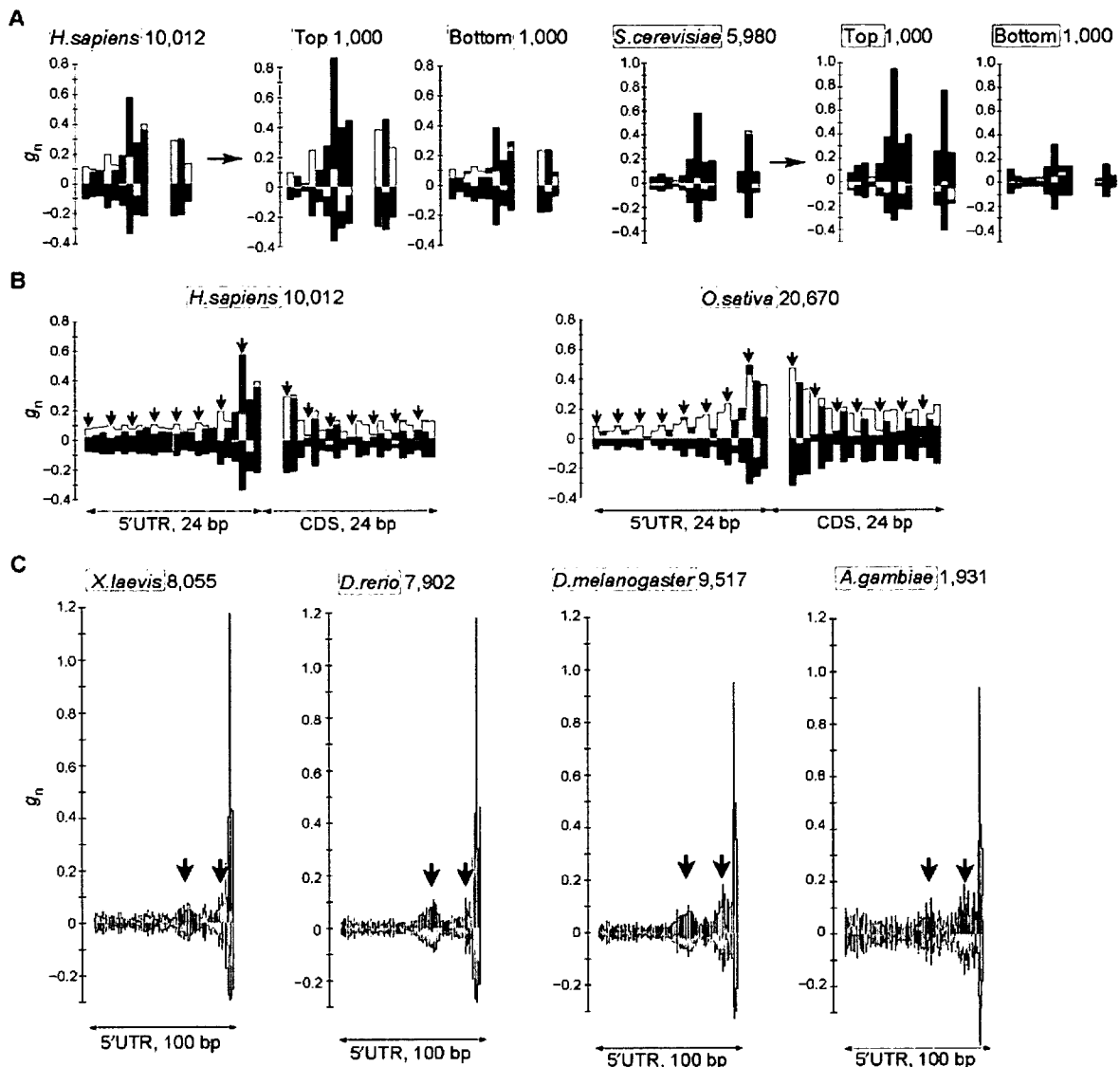


Figure 5. Several features of nucleotide bias around the initiation codon. (A) Nucleotide bias around the initiation codon for genes expressed at high and low levels in humans (left) and *S. cerevisiae* (right). The diagram for each species shown at the left is the same as that in Figure 2. The middle and right diagrams for each species were calculated by using the top 1000 genes with higher expression levels and the bottom 1000 genes with lower expression levels, respectively, in each species. (B) Three-base periodicity observed in humans (left) and *Oryza sativa* (right). Arrows indicate every three bases. (C) U- and A-rich biases observed at positions around -40 and around -15, respectively, which are indicated by arrows. The g_n values are shown from positions -100 to -1 for *Xenopus laevis*, *Danio rerio*, *D. melanogaster* and *Anopheles gambiae*.

Table 2. Fractions (%) of -3A and +5G for genes expressed at high and low levels

	<i>H. sapiens</i>			<i>S. cerevisiae</i>		
	All	Top 1000	Bottom 1000	All	Top 1000	Bottom 1000
-3A	39.3	46.3	36.8	58.2	72.2	46.0
+5C	35.8	41.5	34.4	38.0	50.6	29.2

Table 3. Observed and expected numbers of genes containing a preferred sequence

Pattern	<i>H. sapiens</i>			<i>C. elegans</i>			<i>O. sativa</i>			<i>A. thaliana</i>		
	O	E	Z	O	E	Z	O	E	Z	O	E	Z
GCCGCCAUG	79	16.7	15.3*	4	0.4	5.5*	262	44.6	32.6*	20	1.5	14.9*
AAAAAAAUG	13	2.7	6.3*	62	36.3	4.3*	33	2.8	17.9*	235	73.2	19.0*

O and *E* represent the observed and expected numbers of genes containing a given sequence, respectively. *E* was calculated under the assumption that each nucleotide at each position appears independently. For example, *E* for AAAAAAUG in humans was calculated as $N o_A^{(-6)} o_A^{(-5)} o_A^{(-4)} o_A^{(-3)} o_A^{(-2)} o_A^{(-1)} = 9857 \times 0.215 \times 0.207 \times 0.251 \times 0.393 \times 0.295 \times 0.213 = 2.7$, where *N* is the number of human genes with 5' UTRs that are six or more bases long, and $o_A^{(i)}$ is the observed fraction of A at position *i*. The deviation of *O* from *E* was evaluated by the Z-value (see Materials and Methods section). An asterisk indicates $P < 10^{-4}$.

eukaryotic mRNA directly base pairs with an 18S rRNA to enhance translation initiation, which is similar to the interaction of the SD sequence with a 16S rRNA in a prokaryotic mRNA. It should be noted that the original scanning model cannot account for these observations.

Our results are consistent with the previous assertion that preferred sequences around the initiation codon vary among different eukaryote species (Figure 2) (5,6). However, Figure 4 suggests that the sequences could generally be decomposed into two distinct patterns, the repetition of GCC and that of A. To examine this possibility in more detail, we compared the observed and expected numbers of genes containing the sequences GCCGCCAUG and AAAAAAUG in several species (Table 3). The expected number was calculated based on the assumption that an observed nucleotide at each position will appear in a manner independent of a nucleotide at another position. The results clearly showed that the observed numbers are significantly larger than the expected numbers for these sequences. It is therefore suggested that the sequence GCCGCCAUG or AAAAAAUG, and not a particular nucleotide at each position, may play a role as a whole in translation initiation. We further examined the existence of genes containing a hexanucleotide sequence that is a mixture of these two sequences (e.g. GAAACCAUG or ACAGACAUG). We then found that such mixed sequences are significantly suppressed in genes ($P < 0.01$), while AAAAAAUG and GCCGCCAUG are the most and the second most over-represented patterns, respectively, among all hexanucleotide

Table 4. Observed and expected numbers of genes that contain a preferred sequence or a mixed sequence

Pattern	<i>O</i>	<i>E</i>	<i>Z</i>	Rank
AAAAAAAUG	1725	337.3	75.6	1
GCCGCCAUG	841	113.5	68.3	2
ACAGACAUG	190	142.1	4.0	523
AACACCAUG	242	215.6	1.8	1008
GAAGCAAUG	158	139.6	1.6	1093
GCAACAAUG	274	261.9	0.7	1398
ACAACCAUG	225	240.4	-1.0	2302
ACAGACAUG	146	162.9	-1.3	2473
GCACAAAUG	54	69.1	-1.8	2718
GCAACCAUG	261	300.2	-2.3	2946
GACAACAUG	226	269.4	-2.6	3113
GACACAAUG	185	234.9	-3.3	3345
GAAACCAUG	186	236.3	-3.3	3349
ACCAACAUG	185	239.0	-3.5	3438
GAAGACAUG	115	160.1	-3.6	3465
ACCAACAUG	204	274.0	-4.2	3655
AACGCAAUG	72	127.4	-4.9	3788
AACGACAUG	67	146.1	-6.5	3972
GACGAAAUG	68	159.2	-7.2	4027
ACCGAAAUG	50	162.0	-8.8	4075

These numbers were obtained from 219 496 genes in all of the 47 species examined.

sequences (Table 4). These observations support the idea that there are two distinct patterns of signals for translation initiation.

If we assume the presence of the two distinct patterns of signals, then the variation in preferred sequences among different species could be accounted for by differences in the relative contribution from each pattern. In vertebrates or Monocots, the signal of GCC repeats is relatively strong, whereas in invertebrates or Dicots, the signal of repetition of A is more conspicuous. What, then, determines the relative contribution of each pattern to the preferred sequence in a given species? One factor might be the GC content in the genome. Figure 6A shows the GC content in 5' UTRs and that in the whole genome sequences in 25 species with data for more than 3000 genes. This figure suggests that these species can be classified into two groups, i.e. GC-rich and AT-rich. As shown in Figure 6B, a species belonging to the GC-rich group shows a clear signal of GCC repeats, while an AT-rich species frequently exhibits very strong signals of A. These distinct signals might be recognized by different molecular mechanisms. Kozak (34) herself pointed out that the 'Kozak consensus sequence' is repetitious and that the unit of recognition may be a three-base motif. The three-base periodicity observed in this study might help a ribosome locate the correct reading frame. However, the mechanisms for recognizing the abovementioned signals remain unknown at this stage of research. Additional experimental studies will be required to gain a more precise understanding of the molecular mechanisms of translation initiation in eukaryotes.

SUPPLEMENTARY DATA

Supplementary Data are available at NAR Online.

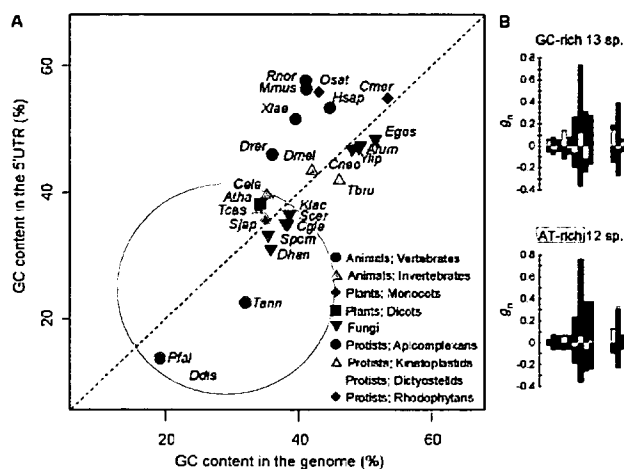


Figure 6. Correlation between GC contents and nucleotide bias around the initiation codon. (A) Horizontal and vertical axes represent the GC contents in the whole genome and in the 5' UTRs, respectively. All species with data for >3000 genes were used. These species can be classified into two groups, i.e. GC-rich (yellow circle) and AT-rich (green circle). We used the genomic GC contents in *X. laevis* (35) and *Schistosoma japonicum* (14) for those in *Xenopus tropicalis* and *Schistosoma mansoni*, respectively, because the genome sequences of *X. laevis* and *S. japonicum* are not available. The name of each species is represented by the initial letter of the generic name and the first three letters of the specific name (Table 1). (B) Bias in nucleotide appearance around the initiation codon for 13 GC-rich species and for 12 AT-rich species. These diagrams were created in the same manner as those in Figure 4.

ACKNOWLEDGEMENTS

We would like to thank Tadashi Imanishi, Motohiko Tanino, Kaoru Mogushi, Takeshi Fukuhara, Emilio Campos, Takeshi Hase, Yutaka Fukuoka, Tadashi Masuda, Soichi Ogishima, and Fengrong Ren for their helpful comments and discussion. Funding for this work was provided by the Ministry of Education, Culture, Sports, Science and Technology of Japan, the Genome Information Integration Project of the Ministry of Economy, Trade and Industry of Japan, and the Japan Biological Informatics Consortium (17710162 to Y.N.). Funding to pay the Open Access publication charges for this article was provided by Tokyo Medical and Dental University.

Conflict of interest statement. None declared.

REFERENCES

- Kozak, M. (1978) How do eucaryotic ribosomes select initiation regions in messenger RNA? *Cell*, **15**, 1109–1123.
- Kozak, M. (2002) Pushing the limits of the scanning mechanism for initiation of translation. *Gene*, **299**, 1–34.
- Cavener, D. and Ray, S. (1991) Eukaryotic start and stop translation sites. *Nucleic Acids Res.*, **19**, 3185–3192.
- Kozak, M. (1984) Compilation and analysis of sequences upstream from the translational start site in eukaryotic mRNAs. *Nucleic Acids Res.*, **12**, 857–872.
- Kozak, M. (1987) An analysis of 5'-noncoding sequences from 699 vertebrate messenger RNAs. *Nucleic Acids Res.*, **15**, 8125–8148.
- Joshi, C., Zhou, H., Huang, X. and Chiang, V. (1997) Context sequences of translation initiation codon in plants. *Plant Mol. Biol.*, **35**, 993–1001.
- Kozak, M. (1986) Point mutations define a sequence flanking the AUG initiator codon that modulates translation by eukaryotic ribosomes. *Cell*, **44**, 283–292.
- Kozak, M. (1999) Initiation of translation in prokaryotes and eukaryotes. *Gene*, **234**, 187–208.
- Cigan, A. and Donahue, T. (1987) Sequence and structural features associated with translational initiator regions in yeast – a review. *Gene*, **59**, 1–18.
- Yun, D., Laz, T., Clements, J. and Sherman, F. (1996) mRNA sequences influencing translation and the selection of AUG initiator codons in the yeast *Saccharomyces cerevisiae*. *Mol. Microbiol.*, **19**, 1225–1239.
- Watanabe, H., Gojobori, T. and Miura, K. (1997) Bacterial features in the genome of *Methanococcus jannaschii* in terms of gene composition and biased base composition in ORFs and their surrounding regions. *Gene*, **205**, 7–18.
- Niimura, Y., Terabe, M., Gojobori, T. and Miura, K. (2003) Comparative analysis of the base biases at the gene terminal portions in seven eukaryote genomes. *Nucleic Acids Res.*, **31**, 5195–5201.
- Imanishi, T., Itoh, T., Suzuki, Y., O'Donovan, C., Fukuchi, S., Koyanagi, K., Barrero, R., Tamura, T., Yamaguchi-Kabata, Y. *et al.* (2004) Integrative annotation of 21,037 human genes validated by full-length cDNA clones. *PLoS Biol.*, **2**, e162.
- Wheeler, D., Barrett, T., Benson, D., Bryant, S., Canese, K., Chetverin, V., Church, D., DiCuccio, M., Edgar, R. *et al.* (2006) Database resources of the National Center for Biotechnology Information. *Nucleic Acids Res.*, **34**, D173–D180.
- Tanino, M., Debily, M., Tamura, T., Hishiki, T., Ogasawara, O., Murakawa, K., Kawamoto, S., Itoh, K., Watanabe, S. *et al.* (2005) The Human Anatomic Gene Expression Library (H-ANGEL), the H-Inv integrative display of human gene expression across disparate technologies and platforms. *Nucleic Acids Res.*, **33**, D567–D572.
- Holstege, F., Jennings, E., Wyrick, J., Lee, T., Hengartner, C., Green, M., Golub, T., Lander, E. and Young, R. (1998) Dissecting the regulatory circuitry of a eukaryotic genome. *Cell*, **95**, 717–728.
- Sokal, R.R. and Rohlf, F.J. (1993) *Biometry*, 3rd edn., 689–697.
- Mewes, H., Frishman, D., Mayer, K., Münsterkötter, M., Noubibou, O., Page, P., Rattei, T., Oesterheld, M., Ruepp, A. *et al.* (2006) MIPS: analysis and annotation of proteins from whole genomes in 2005. *Nucleic Acids Res.*, **34**, D169–D172.
- Hertz-Fowler, C., Peacock, C., Wood, V., Aslett, M., Kerhornou, A., Mooney, P., Tivey, A., Berriman, M., Hall, N. *et al.* (2004) GeneDB: a resource for prokaryotic and eukaryotic organisms. *Nucleic Acids Res.*, **32**, D339–D343.
- Eichinger, L., Pachebat, J., Glöckner, G., Rajandream, M., Sugang, R., Berriman, M., Song, J., Olsen, R., Szafranski, K. *et al.* (2005) The genome of the social amoeba *Dictyostelium discoideum*. *Nature*, **435**, 43–57.
- Lukaszewicz, M., Feuermann, I., Jérôme, B., Stas, A. and Boutry, M. (2000) In vivo evaluation of the context sequence of the translation initiation codon in plants. *Plant Sci.*, **154**, 89–98.
- Kochetov, A. (2005) AUG codons at the beginning of protein coding sequences are frequent in eukaryotic mRNAs with a suboptimal start codon context. *Bioinformatics*, **21**, 837–840.
- Pesole, G., Gissi, C., Grillo, G., Licciulli, F., Liuni, S. and Saccone, C. (2000) Analysis of oligonucleotide AUG start codon context in eukaryotic mRNAs. *Gene*, **261**, 85–91.
- Pisareva, A., Kolupaeva, V., Pisareva, V., Merrick, W., Hellen, C. and Pestova, T. (2006) Specific functional interactions of nucleotides at key -3 and +4 positions flanking the initiation codon with components of the mammalian 48S translation initiation complex. *Genes Dev.*, **20**, 624–636.
- Kozak, M. (1997) Recognition of AUG and alternative initiator codons is augmented by G in position +4 but is not generally affected by the nucleotides in positions +5 and +6. *EMBO J.*, **16**, 2482–2492.

26. Grünert,S. and Jackson,R. (1994) The immediate downstream codon strongly influences the efficiency of utilization of eukaryotic translation initiation codons. *EMBO J.*, **13**, 3618–3630.
27. Sawant,S., Kiran,K., Singh,P. and Tuli,R. (2001) Sequence architecture downstream of the initiator codon enhances gene expression and protein stability in plants. *Plant Physiol.*, **126**, 1630–1636.
28. Morris,D. and Geballe,A. (2000) Upstream open reading frames as regulators of mRNA translation. *Mol. Cell. Biol.*, **20**, 8635–8642.
29. Crowe,M., Wang,X. and Rothnagel,J. (2006) Evidence for conservation and selection of upstream open reading frames suggests probable encoding of bioactive peptides. *BMC Genomics*, **7**, 16.
30. Meijer,H. and Thomas,A. (2002) Control of eukaryotic protein synthesis by upstream open reading frames in the 5'-untranslated region of an mRNA. *Biochem. J.*, **367**, 1–11.
31. Rogers,G.W., Edelman,G.M. and Mauro,V.P. (2004) Differential utilization of upstream AUGs in the beta-secretase mRNA suggests that a shunting mechanism regulates translation. *Proc. Natl Acad. Sci. USA*, **101**, 2794–2799.
32. Lammich,S., Schöbel,S., Zimmer,A., Lichtenthaler,S. and Haass,C. (2004) Expression of the Alzheimer protease BACE1 is suppressed via its 5'-untranslated region. *EMBO Rep.*, **5**, 620–625.
33. Dresios,J., Chappell,S.A., Zhou,W. and Mauro,V.P. (2006) An mRNA-rRNA base-pairing mechanism for translation initiation in eukaryotes. *Nat. Struct. Mol. Biol.*, **13**, 30–34.
34. Kozak,M. (1987) At least six nucleotides preceding the AUG initiator codon enhance translation in mammalian cells. *J. Mol. Biol.*, **196**, 947–950.
35. Birney,E., Andrews,D., Caccamo,M., Chen,Y., Clarke,L., Coates,G., Cox,T., Cunningham,F., Curwen,V. *et al.* (2006) Ensembl 2006. *Nucleic Acids Res.*, **34**, D556–D561.
36. Carninci,P., Kasukawa,T., Katayama,S., Gough,J., Frith,M., Maeda,N., Oyama,R., Ravasi,T., Lenhard,B. *et al.* The transcriptional landscape of the mammalian genome. *Science*, **309**, 1559–1563.
37. Gerhard,D., Wagner,L., Feingold,E., Shenmen,C., Grouse,L., Schuler,G., Klein,S., Old,S., Rasooly,R. *et al.* (2004) The status, quality, and expansion of the NIH full-length cDNA project: the Mammalian Gene Collection (MGC). *Genome Res.*, **14**, 2121–2127.
38. Morin,R., Chang,E., Petrescu,A., Liao,N., Griffith,M., Chow,W., Kirkpatrick,R., Butterfield,Y., Young,A. *et al.* (2006) Sequencing and analysis of 10,967 full-length cDNA clones from *Xenopus laevis* and *Xenopus tropicalis* reveals post-tetraploidization transcriptome remodeling. *Genome Res.*, **16**, 796–803.
39. Rasooly,R., Henken,D., Freeman,N., Tompkins,L., Badman,D., Briggs,J. and Hewitt,A.T. and The National Institutes of Health Trans-NIH Zebrafish Coordinating Committee (2003) Genetic and genomic tools for zebrafish research: the NIH zebrafish initiative. *Dev. Dyn.*, **228**, 490–496.
40. Rhee,S., Beavis,W., Berardini,T., Chen,G., Dixon,D., Doyle,A., Garcia-Hernandez,M., Huala,E., Lander,G. *et al.* (2003) The *Arabidopsis* Information Resource (TAIR): a model organism database providing a centralized, curated gateway to *Arabidopsis* biology, research materials and community. *Nucleic Acids Res.*, **31**, 224–228.
41. Kikuchi,S., Satoh,K., Nagata,T., Kawagashira,N., Doi,K., Kishimoto,N., Yazaki,J., Ishikawa,M., Yamada,H. *et al.* (2003) Collection, mapping, and annotation of over 28,000 cDNA clones from japonica rice. *Science*, **301**, 376–379.
42. Chisholm,R., Gaudet,P., Just,E., Pilcher,K., Fey,P., Merchant,S. and Kibbe,W. (2006) dictyBase, the model organism database for *Dictyostelium discoideum*. *Nucleic Acids Res.*, **34**, D423–D427.
43. Matsuzaki,M., Misumi,O., Shin-i,T., Maruyama,S., Takahara,M., Miyagishima,S., Mori,T., Nishida,K., Yagisawa,F. *et al.* (2004) Genome sequence of the ultrasmall unicellular red alga *Cyanidioschyzon merolae* 10D. *Nature*, **428**, 653–657.

Evola: Ortholog database of all human genes in H-InvDB with manual curation of phylogenetic trees

Akihiro Matsuya^{1,2}, Ryuichi Sakate^{1,3,*}, Yoshihiro Kawahara^{1,3}, Kanako O. Koyanagi⁴, Yoshiharu Sato^{1,3}, Yasuyuki Fujii^{1,3}, Chisato Yamasaki^{1,3}, Takuya Habara^{1,3}, Hajime Nakaoka⁵, Fusano Todokoro^{1,6}, Kaori Yamaguchi^{1,3}, Toshinori Endo⁴, Satoshi Oota⁷, Wojciech Makalowski⁸, Kazuho Ikee⁹, Yoshiyuki Suzuki⁹, Kousuke Hanada⁹, Katsuyuki Hashimoto¹⁰, Momoki Hirai¹¹, Hisakazu Iwama¹², Naruya Saitou¹³, Aiko T. Hiraki^{1,3}, Lihua Jin⁹, Yayoi Kaneko^{1,3}, Masako Kanno^{1,3}, Katsuhiko Murakami^{1,3}, Akiko Ogura Noda^{1,3}, Naomi Saichi^{1,3}, Ryoko Sanbonmatsu^{1,3}, Mami Suzuki^{1,3}, Jun-ichi Takeda^{1,3}, Masayuki Tanaka^{1,3}, Takashi Gojobori^{3,9}, Tadashi Imanishi³ and Takeshi Itoh^{3,14}

¹Integrated Database Group, Japan Biological Information Research Center, Japan Biological Informatics Consortium, ²Government & Public Corporation Information Systems, Hitachi, Co., Ltd., ³Integrated Database Group, Biological Information Research Center, National Institute of Advanced Industrial Science and Technology, Tokyo, ⁴Graduate School of Information Science and Technology, Hokkaido University, Hokkaido, ⁵C's Lab Co., Ltd., Hokkaido, ⁶DYNACOM Co., Ltd, Chiba, ⁷BioResource Center, RIKEN, Ibaraki, Japan, ⁸Institute of Bioinformatics, University of Muenster, Muenster, Germany, ⁹Center for Information Biology and DNA Data Bank of Japan, National Institute of Genetics, Shizuoka, ¹⁰Department of Biomedical Resources, National Institute of Biomedical Innovation, Osaka, ¹¹International Research and Educational Institute for Integrated Medical Sciences, Tokyo Women's Medical University, Tokyo, ¹²Kagawa University, Kagawa, ¹³Department of Population Genetics, National Institute of Genetics, Shizuoka and ¹⁴Division of Genome and Biodiversity Research, National Institute of Agrobiological Sciences, Ibaraki, Japan

Received August 15, 2007; Revised September 27, 2007; Accepted October 1, 2007

ABSTRACT

Orthologs are genes in different species that evolved from a common ancestral gene by speciation. Currently, with the rapid growth of transcriptome data of various species, more reliable orthology information is prerequisite for further studies. However, detection of orthologs could be erroneous if pairwise distance-based methods, such as reciprocal BLAST searches, are utilized. Thus, as a sub-database of H-InvDB, an integrated database of annotated human genes (<http://h-invitational.jp/>), we constructed a fully curated database of evolutionary features of human genes, called 'Evola'. In the process of the ortholog detection, computational analysis based on conserved genome synteny and transcript sequence similarity was followed by manual curation by researchers examining phylogenetic trees. In total, 18968 human genes have orthologs among 11

vertebrates (chimpanzee, mouse, cow, chicken, zebrafish, etc.), either computationally detected or manually curated orthologs. Evola provides amino acid sequence alignments and phylogenetic trees of orthologs and homologs. In 'd_N/d_S view', natural selection on genes can be analyzed between human and other species. In 'Locus maps', all transcript variants and their exon/intron structures can be compared among orthologous gene loci. We expect the Evola to serve as a comprehensive and reliable database to be utilized in comparative analyses for obtaining new knowledge about human genes. Evola is available at <http://www.h-invitational.jp/evola/>.

INTRODUCTION

A large number of genome and transcript sequences accumulated in the last decade give us an opportunity

*To whom correspondence should be addressed. Tel: +81 3 3599 8800; Fax: +81 3 3599 8801; Email: rsakate@ni.aist.go.jp

for large-scale comparative analyses. In particular, detection of orthologs, groups of genes in different species that evolved by speciation, accelerates functional and evolutionary studies. Despite the past efforts to develop bioinformatics methods for analyzing a large number of sequences, it is still a challenge to comprehensively identify orthologs between species. A number of automated pairwise distance-based methods for ortholog detection have been proposed, as represented by the reciprocal best BLAST hits (RBH) method (1) and the reciprocal smallest distance (RSD) method (2). However, as genes might have frequently undergone duplications and losses in evolutionary lineages leading to human (3), pairwise distance-based methods might lead to erroneous inferences of phylogenetic relationships and thus of orthologs. Thus, phylogenetic tree-based detection can be the most plausible solution to provide more reliable orthologs.

Here this database 'Evola', a sub-database complementary to the H-Invitational database (H-InvDB), was developed to provide orthology information for the originally annotated human genes in H-InvDB. Evola features its ortholog detection in which genome synteny-based computational analysis was followed by manual curation of molecular phylogenetic trees. Evola differs in this way from other ortholog databases such as Inparanoid (4), Ensembl-Compara (5), Homologene (6), HOGENOM (7) and TreeFam (8). These databases are based on BLAST hits (Inparanoid), BLAST hits and synteny (Ensembl-Compara and Homologene) and phylogenetic trees (HOGENOM and TreeFam). The concept of Evola is that genomic region (gene locus) is a unit of genes that are duplicated or lost. In collaboration with H-InvDB, Evola enables users to compare gene structure, transcript variants, upstream/downstream region of the genome among species.

H-InvDB is an integrated database of annotated human genes providing annotation of human full-length enriched cDNAs (9,10,11). At the meetings of the Human Full-Length cDNA Annotation Invitational held in Japan (2002 and 2003), Evola started with H-InvDB to annotate evolutionary features of the human genes. With several updates afterwards and a subsequent All Human Genes Evolutionary Annotation (AHG-EV) meeting in 2006, the current strategy of evolutionary annotation (computational analysis and manual curation) in Evola has been established. Orthology information for human and other 11 vertebrates is currently included in the Evola: human, chimpanzee, macaque, mouse, rat, dog, cow, opossum, chicken, zebrafish, Tetraodon and Fugu. Several visualization tools are incorporated into the database, including sequence alignment viewer, natural selection plot and graphical representation of orthologous gene loci among different species. Evola is now one of the databases listed in the Comparison of Orthology Predictions project of the HUGO Gene Nomenclature Committee (HGNC, <http://www.genenames.org/>).

ORTHOLOG DETECTION

Computational analysis: Ortholog detection based on conserved genomic synteny and pairwise distance

Species for ortholog detection were selected with consideration of completeness of their genome assemblies (chromosome level), abundance of transcript sequences (~20 000) and importance in biology (intensively studied or a representative of a phylogenetic clade). Whole genome sequence assemblies of human (hg18), chimpanzee (panTro2), macaque (rheMac2), mouse (mm8), rat (rn4), dog (canFam2), cow (rn4), opossum (monDom4), chicken (galGal3), zebrafish (danRer4), Tetraodon (tetNig1) and Fugu (fr1) were downloaded from UCSC (<http://genome.ucsc.edu/>). Conserved syntenic regions were detected by a modified pairwise genome alignment method (12) using BLASTZ (13) with the options of C = 2, T = 4, Y = 3400 between human and other primates (between more similar genome sequences), and C = 2 between human and non-primate vertebrates (between less similar genome sequences).

For human transcripts, H-InvDB representative transcripts (HITS) were used. Other vertebrates' transcripts (mRNAs) were downloaded from DDBJ (<http://www.ddbj.nig.ac.jp/>) release 66, Ensembl (<http://www.ensembl.org/>) release 38 and RefSeq (<http://www.ncbi.nlm.nih.gov/RefSeq/>) release 17, and their genomic locations (one location per transcript) were detected on cognate genomes by a hybrid method using BLAT (14), BLAST (15) and est2genome (16) as they were used to detect genomic locations of human transcripts in H-InvDB. Representative transcripts (one transcript per gene locus) were determined in consideration of percent identity and coverage to the genome, number of exons, etc. of all transcripts in each locus (9,10,11). Thus, in Evola, representative transcripts were defined as genes.

Lengths of overlapping exons of each gene pair between human and other species were calculated in the genome alignment. A gene pair with the maximum length was selected as the best assignment (not a minimum length was defined). Every gene in a species was assigned to a gene in the other species. If two human genes were assigned to one mouse gene, this was defined as a two-to-one ortholog. As a result, Evola contains not only one-to-one orthologs but also many-to-many orthologs. For all the assignment pairs, coding sequences (CDSs) and amino acid (a.a.) sequences of other species were predicted by FASTY (17). They were predicted by comparing with the amino acid sequences of the corresponding human genes. Finally, if the length of the alignable region between human and other species ortholog candidates was ≥ 80 a.a., they were defined as computationally detected orthologs.

Manual curation: Examination of phylogenetic trees by experts

Homologs of human genes (amino acid sequences) were obtained from UniProt (<http://www.uniprot.org/>) and human RefSeq (NP) by FASTY similarity searches with the option of E-value of $< 1e-5$. For each human gene, a

sequence set consisting of both the computationally detected orthologs and the homologs was prepared. For these sequence sets, phylogenetic trees were constructed by the neighbor-joining (NJ) method (18). In detail, multiple amino acid sequence alignments and phylogenetic trees were constructed by ClustalW (19) with the options of bootstrap = 1000, seed = 1, kimura, tossgaps, bootlabels = node.

Phylogenetic trees were examined by experts in the field of molecular evolution, who attended the evolutionary annotation meetings described in the introduction. The trees were drawn by NJplot (<http://pbil.univ-lyon1.fr/software/njplot.html>) and the default rooting was used. Discarding or re-rooting the tree was judged by the experts if necessary. All the ortholog pairs of human and other species detected by the computational analysis were examined (Figure 1). The primary principles of manual curation in Evola to be checked were as follows. [1] Phylogenetic topology between gene tree and species tree is consistent. As a gene tree, the minimum sub-clade including the pair (a part of the tree) was examined. As a species tree of reference, a phylogenetic tree indicating the trifurcation among primates, rodents and Laurasiatherian (dog, cow, etc.) species (20) was used, because the phylogenetic relationship has been controversial among them (21). In fact, we found that ((human–mouse)–dog) clades for some genes and ((human–dog)–mouse) clades for other genes. [2] Outgroup includes either two or more species that are phylogenetically distant from all the species in the sub-clade, or human and other species. In the latter case, human duplicate genes might exist. [3] Available bootstrap values of the corresponding three branches (one between the sub-clade and outgroup, and its two descendants) are all ≥ 900 . The gene pairs consistent with all the principles were defined as ‘manually curated orthologs’, otherwise their annotation status remained to be ‘computationally detected orthologs’.

DATABASE CONTENTS

Evola contains two ortholog datasets: (1) more comprehensive set of orthologs (computational analysis); and (2) more reliable orthologs (computational analysis supported by manual curation). In the current Evola (release 4.1), orthology information for 18 968 human genes is available among 11 vertebrates: chimpanzee, macaque, mouse, rat, dog, cow, opossum, chicken, zebrafish, Tetraodon and Fugu (Table 1). Manually curated orthologs occupied 25.4% of all computationally detected ortholog pairs (24 122/94 935) (release 4.1, 2007).

Evola is a sub-database of H-InvDB (9,10,11), and orthology information in Evola is, as ‘Evolutionary annotation’, a part of the comprehensive human gene annotations in H-InvDB. Thus, orthology information can be utilized with close reference to other annotation in H-InvDB. For example, 2090 human genes with orthology information belonged to H-Inv protein similarity categories of ‘hypothetical proteins’ (similarity category IV–VI). Molecular functions of these hypothetical

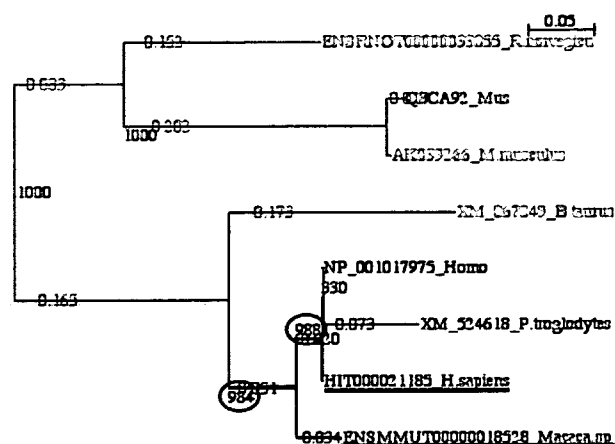


Figure 1. An example of manually curated gene pair from *H.sapiens* (red underline) and *Macaca.sp.* (blue underline). In this case, conditions of phylogenetic topologies, outgroup species (light gray background) and bootstrap values (two circles) are sufficient (refer to the text). Thus, the pair was defined as a manually curated ortholog.

Table 1. Number of orthologs provided in Evola (release 4.1, June 2007)

Species	Genes	Human genes
<i>Homo sapiens</i> (Human)	18 968	–
<i>Pan troglodytes</i> (Chimpanzee)	16 368	15 615
<i>Macaca sp.</i> (Macaque) ^a	12 037	12 352
<i>Mus musculus</i> (Mouse)	15 570	14 574
<i>Rattus norvegicus</i> (Rat)	15 632	14 302
<i>Canis familiaris</i> (Dog)	14 730	13 916
<i>Bos taurus</i> (Cow)	9375	10 181
<i>Monodelphis domestica</i> (Opossum)	13 201	13 588
<i>Gallus gallus</i> (Chicken)	9266	10 738
<i>Danio rerio</i> (Zebrafish)	12 334	10 468
<i>Tetraodon nigroviridis</i> (Tetraodon)	11 505	9820
<i>Takifugu rubripes</i> (Fugu)	9738	9459

Numbers of genes of both human and other species are listed. Owing to lineage-specific duplication or loss, the numbers are usually different (for example, 15 570 mouse genes are orthologous to 14 574 human genes). 18 968 human genes have at least one ortholog among other 11 species.

^a*Macaca mulatta*, *Macaca fascicularis*, *Macaca fuscata*, etc. are included.

proteins can be analyzed using model species. Moreover, cross references between Evola and other annotations in H-InvDB (protein–protein interaction (PPI), expression, polymorphism, disease, etc.) can produce valuable information contributing to the comprehensive understanding of the human genes.

We aimed to develop user-friendly interfaces that provide easy access to a variety of orthology information in Evola. Users can search orthologs in the top page of Evola as well as in the search systems of H-InvDB [simple search, advanced search and navigation system (Navi)]. Users can download data for each human gene on the main page as well as all the data of Evola in the download page. On the main page of Evola (Figure 2),

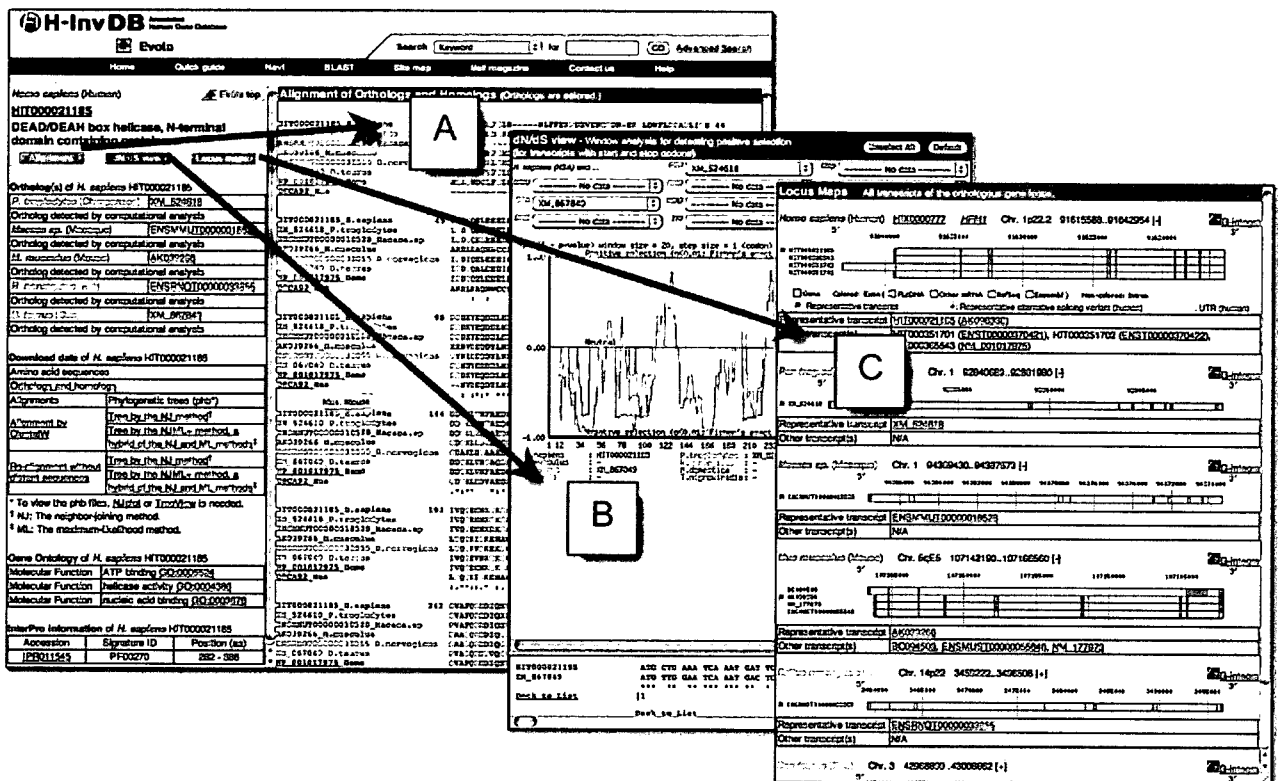


Figure 2. Evola main page. This page is divided into left and right frames. In the left frame, tables of orthologs, download data, Gene ontology, and InterPro are listed. Three green buttons are links to show 'Alignment' (A), ' d_N/d_S view' (B) and 'Locus maps' (C) in the right frame.

the following information for a human gene is available in the left frame: gene name, ortholog list with annotation status, download of sequences, alignments and phylogenetic trees, Gene ontology (22) and InterPro (23). In addition to the set of original ClustalW alignments, another set of alignments, including properly aligned sequences only (24), was also constructed and provided. In the latter sets, sequences with distinctively low identity to other sequences in an alignment were excluded. Based on both alignment sets, phylogenetic trees were constructed by the neighbor-joining method (18) and the NJML+ method (25).

In the right frame of the main page, Evola features the three views described below. Users can switch among the views.

Alignment: Multiple alignments of orthologs and homologs (Figure 2A)

Amino acid sequence alignments of orthologs and homologs are displayed. Users can switch from 'Alignment of Orthologs' (default) to 'Alignment of Orthologs and Homologs', or vice versa. Each amino acid residue is color coded as defined in ClustalX (19). Accession numbers and species names of orthologs (human and other species) are colored in their species colors defined in Evola (human in red, mouse in gray, etc.). Accession numbers of homologs are linked to the

original data sources of UniProt or RefSeq. While species are labeled by their scientific names (Homo, Mus, etc.), users can activate a popup window giving a species common name by placing the mouse cursor over homolog accession numbers (for example, 'Q5R508_Pongo'). InterPro data in the left frame include positional information on a human gene, and they can be utilized to detect conserved domains in the proteins.

d_N/d_S view: Window analysis detecting regions under positive or negative selection (Figure 2B)

Users can select one or more species for which to show the plots in the graph. In the lower frame under the graph, the pairwise nucleotide sequence alignment of CDSs is shown. The sequence positions (a.a. or codon) appearing in the graph and alignment are those of human genes.

The nonsynonymous to synonymous substitution rate ratio (d_N/d_S) is a commonly used measure of natural selection. In order to visualize positively and negatively selected regions, sliding window analysis was conducted (a 20 codon window with 1 codon stepping; result for the first window appears as a plot at 11th codon of the human gene). The statistical significance (P -value) of the difference between the number of nonsynonymous substitution (n) per synonymous substitutions (s); n/s , and the number of nonsynonymous sites (N) per synonymous sites (S); N/S was calculated by Fisher's exact test.

d_S , d_N , s and n values were estimated by the modified Nei–Gojobori method (26,27). If $d_N/d_S > 1$, the score ($= 1 - P$ -value) was plotted above the zero line (neutral), and if $d_N/d_S < 1$, the score [$= -(1 - P$ -value)] was plotted below the zero line. The regions plotted above the red line indicate that the sites might be under positive selection ($d_N/d_S > 1$ and $P < 0.01$). Conversely, the regions plotted below the blue line indicate that the sites might be under negative (purifying) selection ($d_N/d_S < 1$ and $P < 0.01$).

Locus maps: Comparative maps of orthologous gene loci (Figure 2C)

Orthologs were detected for representative transcripts (one transcript per gene locus) in Evola. However, there could be transcript variants in gene loci that have different exon–intron structures leading to produce different protein isoforms. Thus, information on other transcripts besides the representative transcript among orthologous gene loci are shown in Locus maps. In the figures, exon/intron structure, coding sequence (CDS) and untranslated regions (UTR) for each transcript are visualized. H-Inv cluster ID (HIX, an identifier of gene locus), Gene symbol, genomic location and a link to ‘G-integra’, an integrated genome browser of H-InvDB, are available. The flag icon denotes the representative transcript. The blue diamond icon denotes the Representative Alternative Splicing Variant (RASV) that is another representative per transcript group consisting of the same alternative splicing pattern (28). Representative transcripts are also RASVs, and blue diamonds do not appear if there is only one splicing isoform. In the tables, the H-Inv transcript ID (HIX) and original accession numbers (DDBJ/EMBL/GenBank, Ensembl and RefSeq) of the representative transcript and other transcripts are listed.

FUTURE DIRECTIONS

As our update policy, orthology information in Evola is updated when H-InvDB annotation is updated. One major update and three minor updates per year are scheduled. At the next major update on December 2007, a new duplicate gene family view is planned to be integrated within Evola. Human duplicate gene family data was originally constructed based on both amino acid sequence similarity (29) and orthology information. In the current Evola (release 4.1), parts of human duplicate gene annotation have been already implemented. The human duplicate genes are included in the alignments and phylogenetic trees of orthologs and homologs. Finally, we expect Evola to serve as a new database for evolutionary annotation of human genes. We sincerely welcome any requests and feedback from users.

ACKNOWLEDGEMENTS

We thank the members of Integrated Database Group, Japan Biological Information Research Center for their

helpful suggestions. We are also grateful to Craig Gough for critical reading of the manuscript. This work was supported by the Ministry of Economy, Trade and Industry of Japan (METI), and the Japan Biological Informatics Consortium (JBIC). Funding to pay the Open Access publication charges for this article was provided by JBIC.

Conflict of interest statement. None declared.

REFERENCES

- Tatusov, R.L., Koonin, E.V. and Lipman, D.J. (1997) A genomic perspective on protein families. *Science*, **278**, 631.
- Wall, D.P., Fraser, H.B. and Hirsh, A.E. (2003) Detecting putative orthologs. *Bioinformatics*, **19**, 1710–1711.
- Fortna, A., Kim, Y., MacLaren, E., Marshall, K., Hahn, G., Meltesen, L., Brenton, M., Hink, R., Burgers, S., Hernandez-Boussard, T. *et al.* (2004) Lineage-specific gene duplication and loss in human and great ape evolution. *PLoS Biol.*, **2**, E207.
- O'Brien, K.P., Remm, M. and Sonnhammer, E.L. (2005) Inparanoid: a comprehensive database of eukaryotic orthologs. *Nucleic Acids Res.*, **33**, D476–D480.
- Birney, E., Andrews, D., Caccamo, M., Chen, Y., Clarke, L., Coates, G., Cox, T., Cunningham, F., Curwen, V., Cutts, T. *et al.* (2006) Ensembl 2006. *Nucleic Acids Res.*, **34**, D556–D561.
- Wheeler, D.L., Barrett, T., Benson, D.A., Bryant, S.H., Canese, K., Chetverin, V., Church, D.M., DiCuccio, M., Edgar, R., Federhen, S. *et al.* (2007) Database resources of the National Center for Biotechnology Information. *Nucleic Acids Res.*, **35**, D5–D12.
- Dufayard, J.F., Duret, L., Penel, S., Gouy, M., Rechenmann, F. and Perriere, G. (2005) Tree pattern matching in phylogenetic trees: automatic search for orthologs or paralogs in homologous gene sequence databases. *Bioinformatics*, **21**, 2596–2603.
- Li, H., Coghlan, A., Ruan, J., Coin, L.J., Heriche, J.K., Osmotherly, L., Li, R., Liu, T., Zhang, Z., Bolund, L. *et al.* (2006) TreeFam: a curated database of phylogenetic trees of animal gene families. *Nucleic Acids Res.*, **34**, D572–D580.
- Imanishi, T., Itoh, T., Suzuki, Y., O'Donovan, C., Fukuchi, S., Koyanagi, K.O., Barrero, R.A., Tamura, T., Yamaguchi-Kabata, Y., Tanino, M. *et al.* (2004) Integrative annotation of 21,037 human genes validated by full-length cDNA clones. *PLoS Biol.*, **2**, e162.
- Yamasaki, C., Koyanagi, K.O., Fujii, Y., Itoh, T., Barrero, R., Tamura, T., Yamaguchi-Kabata, Y., Tanino, M., Takeda, J., Fukuchi, S. *et al.* (2005) Investigation of protein functions through data-mining on integrated human transcriptome database, H-Invitational database (H-InvDB). *Gene*, **364**, 99–107.
- Yamasaki, C., Murakami, K., Fujii, Y., Sato, Y., Harada, E., Takeda, J., Taniya, T., Sakate, R., Kikugawa, S., Shimada, M. *et al.* (2008) The H-Invitational Database (H-InvDB), A comprehensive annotation resource for human genes and transcripts. *Nucleic Acids Res.*, **36**, in press.
- Fujii, Y., Itoh, T., Sakate, R., Koyanagi, K.O., Matsuya, A., Habara, T., Yamaguchi, K., Kaneko, Y., Gojobori, T., Imanishi, T. *et al.* (2005) A web tool for comparative genomics: G-compass. *Gene*, **364**, 45–52.
- Schwartz, S., Kent, W.J., Smit, A., Zhang, Z., Baertsch, R., Hardison, R.C., Haussler, D. and Miller, W. (2003) Human-mouse alignments with BLASTZ. *Genome Res.*, **13**, 103–107.
- Kent, W.J. (2002) BLAT—the BLAST-like alignment tool. *Genome Res.*, **12**, 656–664.
- Altschul, S.F., Gish, W., Miller, W., Myers, E.W. and Lipman, D.J. (1990) Basic local alignment search tool. *J. Mol. Biol.*, **215**, 403–410.
- Rice, P., Longden, I. and Bleasby, A. (2000) EMBOS: the European Molecular Biology Open Software Suite. *Trends Genet.*, **16**, 276–277.
- Pearson, W.R. (2000) Flexible sequence similarity searching with the FASTA3 program package. *Methods Mol. Biol.*, **132**, 185–219.