

from the pillar fell in the direction of *Anantaśāyīn* Vishnu in cave 13. Moreover, on the important summer solstice day at Udayagiri, the motion of the sun across the sky was in line with the specially cut passageway and the evening setting sun completely illuminated the *Anantaśāyīn* Vishnu panel in cave 13. The image that was probably atop the Delhi iron pillar capital has been deduced to be disc-shaped, approximately 20" in diameter and 2" thick. The nature of the top surface of the iron pillar capital indicated that the disc-shaped object was fitted vertically on a flat, circular base, which was subsequently welded onto the top of the cylinder, around which the components of the decorative bell capital were shrunk fit.

1. Anantharaman, T. R., *The Rustless Wonder – A Study of the Delhi Iron Pillar*, Vigyan Prasar, New Delhi, India, 1997.
2. Balasubramaniam, R., *Delhi Iron Pillar: New Insights*, Indian Institute of Advanced Study, Shimla, 2002, pp. 8–46.
3. Balasubramaniam, R., Identity of *Chandra* and *Vishnupadagiri* of the Delhi iron pillar inscription: numismatic, archaeological and literary evidence. *Bull. Met. Mus.*, 2000, **32**, 42–64.
4. Dass, M. I., *Udayagiri – Rise of a sacred hill, its art, architecture and landscape – A study*. Ph D thesis, DeMontfort University, Leicestershire, UK, 2001, pp. 105–155.
5. Dass, M. I. and Balasubramaniam, R., Estimation of the original erection site of the Delhi iron pillar at Udayagiri. *Indian J. Hist. Sci.*, 2004, **39.1**, 51–74.
6. Willis, M., Inscriptions at Udayagiri: Locating domains of devotion, patronage and power in the eleventh century. *South Asian Stud.*, 2001, **17**, 48.
7. Lahiri, N., *Archaeology of Ancient Indian Trade Routes upto c. 300 BC*, Oxford University Press, New Delhi, India, 1992, pp. 368–381.
8. Dass, M. I. and Willis, M., The lion capital from Udayagiri and the antiquity of sun worship in central India. *South Asian Stud.*, 2002, **18**, 25–45.
9. Gangooly, P. (ed.), *Surya Siddhanta, Translated with Notes and Appendix by Rev. Ebenezer Burgess (1880)*, Motilal Banarsidass, New Delhi, 1960.
10. Thibaut, G., Contribution to the explanation of the *Jyotisa-Vedanga*. *J. Asiat. Soc. Bengal*, 1977, **46**, 411–437.
11. Krishna Deva, Gupta and their feudatories. In *Encyclopaedia of Indian Temple Architecture*, North India: Foundations of North Indian Style c 250 BC–1100 AD (eds Meister, M. W., Dhaky M. A. and Krishna Deva), American Institute of Indian Studies, Princeton University Press, 1988, vol. II, part I, pp. 18–57.
12. Patil, D. R., *The Monuments of the Udayagiri Hill*, Gwalior, 1948, pp. 377–428.
13. Mitra, D., Varāha cave at Udayagiri – an iconographic study. *J. Asiat. Soc.*, 1963, **5**, 99–103.
14. Williams, J., *The Art of Gupta India: Empire and Province*, Princeton University Press, Princeton, USA, 1982, pp. 37–49.
15. Fleet, J. F., Inscriptions of the early Gupta kings and their successors. *Corpus Inscriptionum Indicarum*, 1888, vol. III, p. 25.
16. Majumdar, R. C., In *The History and Culture of the Indian People – Classical Age*, Bharatiya Vidya Bhavan, Mumbai, 1954, vol. 3, pp. 321–323.
17. Balasubramaniam, R., Dass, M. I. and Raven, E. M., On the original image atop the Delhi iron pillar. *Indian J. Hist. Sci.*, 2004, **39.2**, in press.
18. Balasubramaniam, R., Decorative bell capital of the Delhi iron pillar. *J. Met.*, 1998, **50**, 40–47.

19. Balasubramaniam, R., New insights on the corrosion of the Delhi iron pillar based on historical and dimensional analysis. *Curr. Sci.*, 1997, **73**, 1057–1067.
20. Agrawal, V. S., *Chakra-dhwaja: The Wheel Flag of India*, Prithvi Prakashan, Varanasi, 1964, pp. 1–59.

ACKNOWLEDGEMENTS. We thank the Archaeological Survey of India for co-operation and Dr Ellen M. Raven for critical discussions.

Received 5 April 2003; revised accepted 15 December 2003

Shannon's uncertainty principle and gene expression levels

V. Subramaniam[†], S. K. Gupta and T. C. Ghosh*

Bioinformatics Centre, Bose Institute, P/12, CIT Scheme VII M, Kolkata 700 054, India

[†]Present address: Indian Institute of Information Technology, Allahabad 211 012, India

Shannon's uncertainty principle has been applied to measure the degree of constraints in codon bias in the coding sequences of *Escherichia coli*, *Saccharomyces cerevisiae* and *Haemophilus influenzae*. Our study shows a high degree of correlation between Shannon's uncertainty values and codon adaptation index (CAI). This result suggests that the degree of constraints determined by Shannon's uncertainty principle reflects the level of gene expression and we propose that Shannon's uncertainty principle can be used as an alternative method for predicting the level of gene expression. The main advantage of Shannon's uncertainty values over CAI is that for calculating Shannon's uncertainty values, one does not require any reference set of genes as required in CAI. This gives a potential use of Shannon's uncertainty principle over CAI in predicting the level of gene expression, specially for the newly sequenced genomes where genes are not properly annotated.

THE existence of non-random uses of synonymous codons is well documented. Moreover, codon usage patterns differ significantly among different genes within the same taxon¹. It has been widely accepted that compositional biases are the only dictators in shaping the codon usage variation among the genes in the extremely AT or GC-rich unicellular organisms^{2–4}. It has been suggested that translational selection determines the codon usage bias of highly expressed genes and subsequently, it has been advocated that preferred codons in highly expressed genes are recognized by most abundant tRNAs^{5–7}. Recently, it was

*For correspondence. (e-mail: tapash@bic.boseinst.ernet.in)

observed that in *Pseudomonas aeruginosa*, codon usage bias is mainly dictated by translational selection rather than mutational biases, though it is an extremely GC-rich organism⁸. In some unicellular organisms, it was observed that both translational and compositional constraints are operational in dictating the codon usage variation among the genes in those organisms^{5,9–11}. In *Borrelia burgdorferi*, it was observed that replicational–transcriptional selection is responsible for the codon usage variation among the genes¹². Recently, it has been reported that the cellular as well as the physical location of the gene products can also reflect the codon usage patterns^{13,14}. In *Mycobacterium*, it has been reported that codon usage bias was mainly dictated by hydrophobicity of each gene¹⁵. Codon usage was found to be affected by the base composition of the neighbouring sites¹⁶. Possible relationships between synonymous codon usage and protein secondary structures have also been reported in the literature^{17–20}.

There may be other unknown constraints coming into play in selecting the codon usage variation in different organisms and availability of large number of complete genomes opens up a tremendous opportunity to reveal those unknown factors. Several mathematical indices have been proposed to estimate the degree of codon bias and until now, the best method to capture gene expression using translational bias is codon adaptation index (CAI), as described by Sharp and Li²¹. CAI takes into consideration that the highly expressed genes use optimal codons more frequently than other genes in the genome and computes a weight for each codon and combines them to define the CAI value for each gene in the genome.

Although CAI has been widely used in predicting highly expressed genes, the interpretation of the results requires much more care. In species where translational selection is either absent or ineffective and mutational bias is more pronounced, the CAI values for individual genes are related to its base composition rather than its expression level. The other drawback of CAI is that it requires prior knowledge of optimal codons or a reference set of highly expressed genes in the species.

It has become necessary to devise new methods, so that one can predict the level of gene expression without any prior knowledge of optimal codons in the highly expressed genes. With this aim, we have used Shannon's uncertainty principle to predict the level of gene expression. Shannon's uncertainty value in bits per codon, is used to know the degree of constraint present in the usage of codons while coding the sequence^{22–24}.

The complete genomes of *Escherichia coli*, *Saccharomyces cerevisiae* and *Haemophilus influenzae* have been downloaded from ftp.ncbi.nlm.nih.gov/genbank/genomes. Our own program developed in C was used to retrieve the coding sequences from the complete genomes. To minimize sampling errors, we have chosen only those sequences that are greater than or equal to 300 bp and have correct initial and termination codons. Finally, 3892, 4725

and 1568 genes were selected for *E. coli*, *S. cerevisiae* and *H. influenzae* respectively.

Shannon's uncertainty values (H) were computed using

$$H = \left(- \sum_{i=1}^{i=61} p_i \log_2 P_i \right), \quad (1)$$

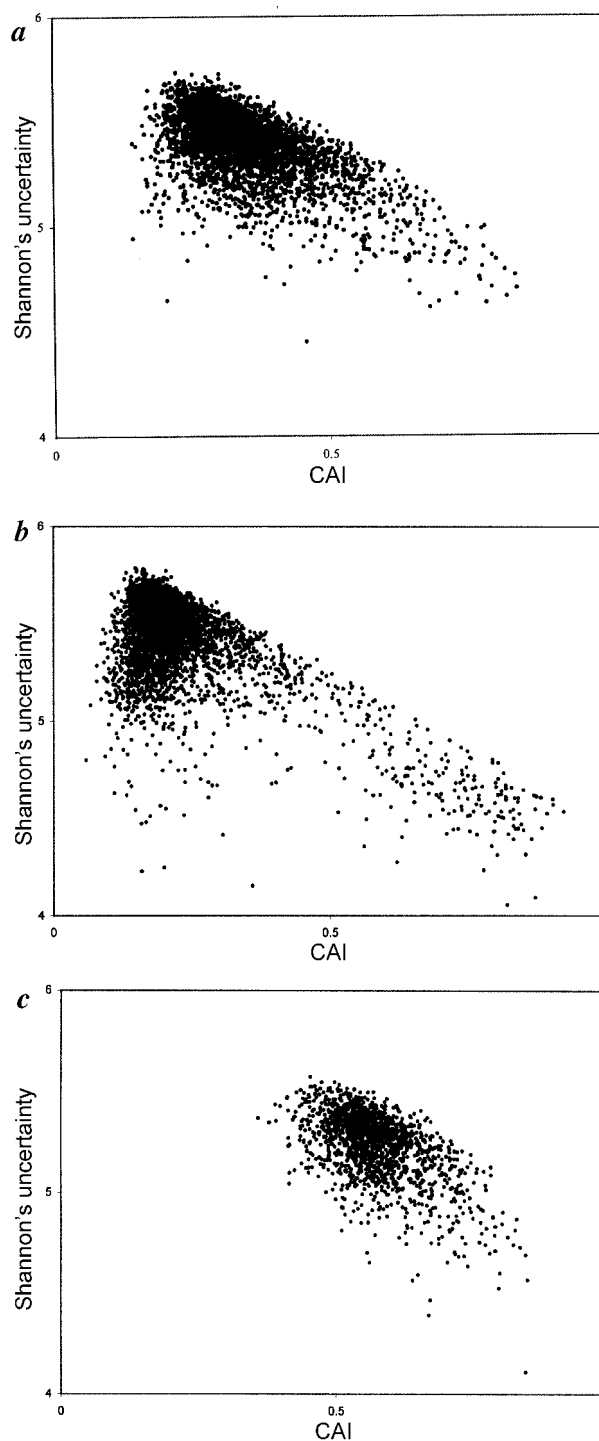


Figure 1. Scatter plot between uncertainty and CAI values of (a) *E. coli*, (b) *S. cerevisiae* and (c) *H. influenzae* genes.

RESEARCH COMMUNICATIONS

where p_i is the probability of a codon present in the given gene and i runs from 1 to 61. When a codon is found to be absent for a gene, we introduced a Laplace correction of 1.0×10^{-6} for p_i . We have excluded three termination codons in our analysis. H gives the uncertainty or degree of constraint in bits per codon present in a given gene in accordance with their codon usage²². In a gene, when all

the codons are used in equal probability then according to the eq. (1), the uncertainty value will be 5.93 bits per codon and any value less than 5.93 bits per codon implies that certain degrees of constraint are present in that gene.

CAI was calculated for all the genes in the study according to the method adopted by Sharp and Li²¹. The reference subset to calculate CAI for *E. coli* and *S. cere-*

Table 1. Uncertainty and CAI values of 50 genes located at the extreme end in the uncertainty scale of lowest uncertainty values

| <i>E. coli</i> | | | <i>S. cerevisiae</i> | | | <i>H. influenzae</i> | | |
|----------------|-------------------------------|------|----------------------|-------------------------------|------|--|-------------------------------|------|
| Gene | Uncertainty in bits per codon | CAI | Gene | Uncertainty in bits per codon | CAI | Gene | Uncertainty in bits per codon | CAI |
| <i>rplL</i> | 3.99 | 0.85 | <i>RPP2A</i> | 4.06 | 0.82 | <i>rpL7/L12</i> | 4.11 | 0.85 |
| <i>asr</i> | 4.45 | 0.46 | <i>CCW12</i> | 4.09 | 0.88 | Outer membrane protein | 4.53 | 0.80 |
| <i>rplT</i> | 4.61 | 0.69 | <i>RPL42B</i> | 4.23 | 0.78 | <i>rpL9</i> | 4.57 | 0.85 |
| <i>rpsI</i> | 4.63 | 0.79 | <i>RPL36B</i> | 4.31 | 0.81 | <i>rpL22</i> | 4.60 | 0.80 |
| <i>rplO</i> | 4.64 | 0.70 | <i>RPS12</i> | 4.31 | 0.86 | <i>rpL15</i> | 4.63 | 0.74 |
| <i>ompC</i> | 4.66 | 0.83 | <i>FIT3</i> | 4.36 | 0.56 | <i>rpL20</i> | 4.68 | 0.73 |
| <i>rplI</i> | 4.68 | 0.73 | <i>RPS26A</i> | 4.36 | 0.80 | <i>rpL21</i> | 4.68 | 0.74 |
| <i>rplW</i> | 4.68 | 0.67 | <i>RPS20</i> | 4.39 | 0.84 | <i>rpL1</i> | 4.69 | 0.85 |
| <i>eno</i> | 4.70 | 0.84 | <i>RPL30</i> | 4.39 | 0.87 | <i>bolA</i> | 4.69 | 0.72 |
| <i>mopA</i> | 4.71 | 0.80 | <i>TIR1</i> | 4.40 | 0.63 | <i>rrf</i> | 4.70 | 0.78 |
| <i>glnB</i> | 4.72 | 0.42 | <i>RPL24B</i> | 4.42 | 0.77 | <i>atpF</i> | 4.70 | 0.72 |
| <i>tsf</i> | 4.74 | 0.78 | <i>RPS24A</i> | 4.42 | 0.83 | <i>rpL19</i> | 4.71 | 0.79 |
| <i>rplA</i> | 4.76 | 0.78 | <i>RPS26B</i> | 4.42 | 0.71 | <i>pal</i> | 4.73 | 0.84 |
| <i>gapA</i> | 4.77 | 0.84 | <i>RPS23B</i> | 4.43 | 0.80 | <i>tsf</i> | 4.74 | 0.83 |
| <i>tufA</i> | 4.79 | 0.82 | <i>RPL24A</i> | 4.43 | 0.76 | <i>rpL11</i> | 4.74 | 0.78 |
| <i>rpsA</i> | 4.81 | 0.78 | <i>RPS10B</i> | 4.43 | 0.80 | <i>hns</i> | 4.75 | 0.69 |
| <i>tig</i> | 4.82 | 0.74 | <i>RPS6B</i> | 4.44 | 0.85 | <i>secG</i> | 4.75 | 0.68 |
| <i>yfiD</i> | 4.82 | 0.71 | <i>RPS25A</i> | 4.44 | 0.73 | <i>iscU</i> | 4.75 | 0.76 |
| <i>slyD</i> | 4.83 | 0.69 | <i>HHF1</i> | 4.45 | 0.74 | <i>rpL24</i> | 4.77 | 0.69 |
| <i>ompX</i> | 4.84 | 0.75 | <i>RPS10A</i> | 4.45 | 0.84 | <i>rplQ</i> | 4.78 | 0.73 |
| <i>rplJ</i> | 4.84 | 0.65 | <i>RPS6A</i> | 4.45 | 0.84 | <i>hupA</i> | 4.78 | 0.70 |
| <i>ahpC</i> | 4.84 | 0.81 | <i>RPL22A</i> | 4.45 | 0.89 | <i>rpL13</i> | 4.79 | 0.79 |
| <i>hlpA</i> | 4.84 | 0.64 | <i>RPS9B</i> | 4.46 | 0.84 | <i>rpS9</i> | 4.80 | 0.77 |
| <i>rplU</i> | 4.85 | 0.65 | <i>RPL8B</i> | 4.47 | 0.85 | <i>secE</i> | 4.80 | 0.63 |
| <i>rplX</i> | 4.85 | 0.59 | <i>RPS14A</i> | 4.48 | 0.80 | <i>rsgA</i> | 4.81 | 0.71 |
| <i>rplK</i> | 4.86 | 0.72 | <i>HSP12</i> | 4.49 | 0.64 | <i>rpS1</i> | 4.81 | 0.83 |
| <i>tufB</i> | 4.86 | 0.80 | <i>TIR2</i> | 4.49 | 0.57 | <i>rps12</i> | 4.82 | 0.72 |
| <i>rplS</i> | 4.86 | 0.65 | <i>RPL8A</i> | 4.50 | 0.84 | <i>trxM</i> | 4.82 | 0.78 |
| <i>rpsF</i> | 4.86 | 0.69 | <i>RPL18A</i> | 4.50 | 0.82 | <i>rpL18</i> | 4.83 | 0.69 |
| <i>rpsJ</i> | 4.86 | 0.58 | <i>RPL17A</i> | 4.50 | 0.81 | <i>rpL3</i> | 4.83 | 0.79 |
| <i>gcvH</i> | 4.86 | 0.55 | <i>RPL13B</i> | 4.51 | 0.75 | <i>mscL</i> | 4.84 | 0.63 |
| <i>yaiU</i> | 4.87 | 0.72 | <i>RPS17A</i> | 4.51 | 0.81 | <i>eno</i> | 4.84 | 0.81 |
| <i>spy</i> | 4.87 | 0.66 | <i>RPL36A</i> | 4.51 | 0.68 | <i>nusG</i> | 4.85 | 0.72 |
| <i>rpsL</i> | 4.87 | 0.67 | <i>RPL25</i> | 4.52 | 0.76 | <i>rpS7</i> | 4.85 | 0.77 |
| <i>crr</i> | 4.87 | 0.62 | <i>SRP40</i> | 4.52 | 0.24 | <i>recX</i> | 4.85 | 0.66 |
| <i>rplP</i> | 4.87 | 0.62 | <i>RPL19B</i> | 4.52 | 0.72 | <i>greA</i> | 4.86 | 0.67 |
| <i>rplC</i> | 4.87 | 0.72 | <i>RPL1B</i> | 4.52 | 0.85 | <i>rpL10</i> | 4.86 | 0.73 |
| <i>frdD</i> | 4.88 | 0.51 | <i>RPL34B</i> | 4.52 | 0.79 | <i>fic</i> | 4.86 | 0.72 |
| <i>hns</i> | 4.88 | 0.60 | <i>ENO2</i> | 4.52 | 0.90 | <i>gapdH</i> | 4.87 | 0.83 |
| <i>csgA</i> | 4.89 | 0.57 | <i>TIR3</i> | 4.53 | 0.52 | <i>rpS11</i> | 4.87 | 0.75 |
| <i>tpiA</i> | 4.89 | 0.75 | <i>TDH3</i> | 4.53 | 0.93 | <i>glnB</i> | 4.88 | 0.58 |
| <i>pal</i> | 4.89 | 0.69 | <i>RPS24B</i> | 4.53 | 0.76 | <i>hlpA</i> | 4.88 | 0.79 |
| <i>tolA</i> | 4.89 | 0.40 | <i>RPL32</i> | 4.53 | 0.83 | Opacity protein | 4.88 | 0.63 |
| <i>rpsC</i> | 4.89 | 0.74 | <i>RPS31</i> | 4.53 | 0.81 | <i>tig</i> | 4.88 | 0.79 |
| <i>rpsB</i> | 4.90 | 0.78 | <i>RPL15A</i> | 4.53 | 0.79 | <i>tolR</i> | 4.88 | 0.60 |
| <i>rplV</i> | 4.90 | 0.56 | <i>RPS23A</i> | 4.54 | 0.73 | <i>FkbP</i> -type peptidyl-prolyl cis-trans isomerase | 4.88 | 0.69 |
| <i>nirD</i> | 4.90 | 0.49 | <i>STM1</i> | 4.54 | 0.74 | <i>slyD</i> | 4.89 | 0.72 |
| <i>rplM</i> | 4.90 | 0.68 | <i>RPL19A</i> | 4.54 | 0.71 | <i>rpS14</i> | 4.89 | 0.68 |
| <i>adk</i> | 4.91 | 0.65 | <i>RPL1A</i> | 4.54 | 0.84 | <i>rpS5</i> | 4.90 | 0.73 |
| <i>pgk</i> | 4.92 | 0.73 | <i>ENO1</i> | 4.55 | 0.88 | <i>exbD</i> | 4.90 | 0.63 |

Table 2. Uncertainty and CAI values of 50 genes located at the extreme end in the uncertainty scale of highest uncertainty values

| <i>E. coli</i> | | | <i>S. cerevisiae</i> | | | <i>H. influenzae</i> | | |
|----------------|-------------|------|----------------------|-------------|------|--|-------------|------|
| Gene | Uncertainty | CAI | Gene | Uncertainty | CAI | Gene | Uncertainty | CAI |
| <i>pqiA</i> | 5.60 | 0.27 | <i>MAL13</i> | 5.69 | 0.17 | <i>hrpa</i> | 5.47 | 0.51 |
| <i>idnD</i> | 5.60 | 0.31 | <i>BUB2</i> | 5.69 | 0.16 | Predicted serine protease | 5.47 | 0.54 |
| <i>ilvI</i> | 5.61 | 0.30 | <i>MIP1</i> | 5.69 | 0.17 | <i>sdaA</i> | 5.47 | 0.55 |
| <i>plsX</i> | 5.61 | 0.23 | <i>RIB2</i> | 5.69 | 0.20 | <i>comM</i> | 5.47 | 0.47 |
| <i>sgcR</i> | 5.61 | 0.29 | <i>MHP1</i> | 5.69 | 0.21 | <i>mepA</i> | 5.47 | 0.52 |
| <i>yicK</i> | 5.61 | 0.25 | <i>BNA2</i> | 5.69 | 0.19 | Heme-hemopexin utilization protein C | 5.47 | 0.54 |
| <i>wecB</i> | 5.61 | 0.32 | <i>AQR2</i> | 5.69 | 0.19 | Biotin synthesis protein | 5.47 | 0.45 |
| <i>rseB</i> | 5.61 | 0.29 | <i>BOP1</i> | 5.69 | 0.21 | Integrase/recombinase | 5.47 | 0.41 |
| <i>NadB</i> | 5.61 | 0.31 | <i>AKR2</i> | 5.70 | 0.17 | <i>rbsR</i> | 5.47 | 0.48 |
| <i>FucK</i> | 5.61 | 0.30 | <i>RGT1</i> | 5.70 | 0.17 | <i>nifR3</i> | 5.48 | 0.53 |
| <i>PyrC</i> | 5.62 | 0.37 | <i>AYT1</i> | 5.70 | 0.20 | <i>topB</i> | 5.48 | 0.51 |
| <i>FdhD</i> | 5.62 | 0.29 | <i>TAH11</i> | 5.70 | 0.15 | <i>dnaE</i> | 5.48 | 0.53 |
| <i>PhrB</i> | 5.62 | 0.29 | <i>THI2</i> | 5.70 | 0.17 | Conserved hypothetical protein | 5.48 | 0.45 |
| <i>DacB</i> | 5.62 | 0.28 | <i>GPI16</i> | 5.70 | 0.21 | <i>truB</i> | 5.49 | 0.47 |
| <i>rpiR</i> | 5.62 | 0.27 | <i>UBP7</i> | 5.70 | 0.19 | <i>rec2</i> | 5.49 | 0.46 |
| <i>MalY</i> | 5.62 | 0.35 | <i>SNG1</i> | 5.70 | 0.19 | <i>emrB</i> | 5.49 | 0.52 |
| <i>intA</i> | 5.62 | 0.30 | <i>FRE6</i> | 5.70 | 0.17 | <i>spoT</i> | 5.49 | 0.53 |
| <i>MolR_1</i> | 5.62 | 0.27 | <i>BUD7</i> | 5.70 | 0.18 | Conserved hypothetical transmembrane protein | 5.49 | 0.53 |
| <i>YdaU</i> | 5.62 | 0.23 | <i>DPH2</i> | 5.70 | 0.18 | <i>modC</i> | 5.49 | 0.46 |
| <i>iap</i> | 5.62 | 0.26 | <i>PTK1</i> | 5.70 | 0.18 | Predicted Zn-dependent protease | 5.49 | 0.52 |
| <i>hyfI</i> | 5.62 | 0.26 | <i>DAL1</i> | 5.70 | 0.19 | <i>rffG</i> | 5.49 | 0.52 |
| <i>MolR_2</i> | 5.62 | 0.25 | <i>MMP1</i> | 5.71 | 0.19 | <i>HI1500</i> | 5.49 | 0.44 |
| <i>rcsC</i> | 5.62 | 0.29 | <i>CDC14</i> | 5.71 | 0.18 | <i>hisB</i> | 5.49 | 0.48 |
| <i>AcnA</i> | 5.62 | 0.34 | <i>HST3</i> | 5.71 | 0.19 | <i>trmD</i> | 5.49 | 0.43 |
| <i>ThiH</i> | 5.63 | 0.30 | <i>CRR1</i> | 5.71 | 0.17 | <i>menC</i> | 5.49 | 0.52 |
| <i>rfaQ</i> | 5.63 | 0.26 | <i>UBP5</i> | 5.71 | 0.19 | <i>argH</i> | 5.50 | 0.56 |
| <i>ampC</i> | 5.63 | 0.29 | <i>PTR3</i> | 5.71 | 0.18 | <i>cyaA</i> | 5.50 | 0.55 |
| <i>intE</i> | 5.63 | 0.24 | <i>SGA1</i> | 5.71 | 0.18 | <i>HI0129</i> | 5.50 | 0.49 |
| <i>PabC</i> | 5.63 | 0.25 | <i>TOP3</i> | 5.72 | 0.18 | <i>purL</i> | 5.50 | 0.57 |
| <i>DnaG</i> | 5.64 | 0.29 | <i>ISC1</i> | 5.72 | 0.20 | <i>nifS</i> protein | 5.50 | 0.49 |
| <i>smf</i> | 5.64 | 0.20 | <i>MET30</i> | 5.72 | 0.20 | <i>ung</i> | 5.50 | 0.51 |
| <i>nac</i> | 5.64 | 0.24 | <i>ZMS1</i> | 5.72 | 0.19 | <i>HI1410</i> | 5.50 | 0.49 |
| <i>fimI</i> | 5.64 | 0.21 | <i>SMF2</i> | 5.72 | 0.15 | <i>hsdR</i> | 5.50 | 0.47 |
| <i>FimD</i> | 5.64 | 0.25 | <i>NDC1</i> | 5.72 | 0.17 | <i>glgA</i> | 5.51 | 0.51 |
| <i>HofB</i> | 5.65 | 0.27 | <i>GAL4</i> | 5.72 | 0.19 | <i>muL</i> | 5.51 | 0.47 |
| <i>fecE</i> | 5.65 | 0.23 | <i>LSB6</i> | 5.73 | 0.18 | Integrase | 5.51 | 0.44 |
| <i>intB</i> | 5.65 | 0.24 | <i>SYG1</i> | 5.73 | 0.21 | Conserved hypothetical protein | 5.51 | 0.49 |
| <i>amn</i> | 5.66 | 0.30 | <i>CAC2</i> | 5.73 | 0.13 | <i>glgX</i> | 5.51 | 0.54 |
| <i>lhr</i> | 5.66 | 0.26 | <i>PDR10</i> | 5.73 | 0.19 | <i>glgB</i> | 5.51 | 0.53 |
| <i>AppA</i> | 5.66 | 0.27 | <i>PHM7</i> | 5.73 | 0.21 | Potassium/copper-transporting ATPase | 5.51 | 0.50 |
| <i>umuC</i> | 5.66 | 0.25 | <i>RAX1</i> | 5.73 | 0.16 | Conserved hypothetical protein | 5.52 | 0.52 |
| <i>intF</i> | 5.66 | 0.26 | <i>GFD2</i> | 5.73 | 0.16 | <i>dprA</i> | 5.52 | 0.49 |
| <i>DsdC</i> | 5.66 | 0.27 | <i>STB5</i> | 5.74 | 0.20 | Conserved hypothetical protein | 5.52 | 0.47 |
| <i>CadC</i> | 5.66 | 0.27 | <i>GUT1</i> | 5.74 | 0.22 | <i>priA</i> | 5.52 | 0.49 |
| <i>PhoH</i> | 5.66 | 0.27 | <i>TDP1</i> | 5.74 | 0.17 | <i>bioF</i> | 5.52 | 0.46 |
| <i>fes</i> | 5.67 | 0.28 | <i>FMT1</i> | 5.74 | 0.14 | Conserved hypothetical protein | 5.52 | 0.45 |
| <i>EvgS</i> | 5.67 | 0.24 | <i>CVT17</i> | 5.74 | 0.18 | Transport protein | 5.53 | 0.46 |
| <i>tra8_1</i> | 5.68 | 0.21 | <i>RDR1</i> | 5.76 | 0.16 | Conserved hypothetical protein | 5.53 | 0.51 |
| <i>HipA</i> | 5.68 | 0.25 | <i>PUS2</i> | 5.76 | 0.13 | <i>hisC</i> | 5.55 | 0.48 |
| <i>SgcX</i> | 5.69 | 0.29 | <i>TRK2</i> | 5.77 | 0.21 | <i>HI1522</i> | 5.55 | 0.50 |

visiae was taken from Sharp and Li²¹ and Velculescu *et al.*²⁵ respectively. Due to non-availability of information regarding experimentally known highly expressed genes in *H. influenzae*, we have used all the ribosomal proteins present in the genome as a reference subset to calculate CAI.

Shannon's uncertainty values in bits per codon as well as their corresponding CAI values were calculated for *E. coli*, *S. cerevisiae* and *H. influenzae* genes. Figure 1 a–c shows the scattered plots between uncertainty and CAI values for the genes of *E. coli*, *S. cerevisiae* and *H. influenzae* respectively. Pearson correlation coefficients at 0.01

Table 3. Uncertainty and CAI values for ribosomal proteins of *E. coli*, *S. cerevisiae* and *H. influenzae*

| <i>S. cerevisiae</i> | | | <i>E. coli</i> | | | <i>H. influenzae</i> | | |
|----------------------|-------------------------------|------|----------------|-------------------------------|------|----------------------|-------------------------------|------|
| Gene | Uncertainty in bits per codon | CAI | Gene | Uncertainty in bits per codon | CAI | Gene | Uncertainty in bits per codon | CAI |
| <i>RPP2A</i> | 4.06 | 0.82 | <i>rplL</i> | 3.99 | 0.85 | <i>rpL7/L12</i> | 4.11 | 0.85 |
| <i>RPL42B</i> | 4.23 | 0.78 | <i>rplT</i> | 4.61 | 0.69 | <i>rpL9</i> | 4.57 | 0.85 |
| <i>RPL36B</i> | 4.31 | 0.81 | <i>rpsI</i> | 4.63 | 0.79 | <i>rpL22</i> | 4.60 | 0.80 |
| <i>RPS12</i> | 4.31 | 0.86 | <i>rplO</i> | 4.64 | 0.70 | <i>rpL15</i> | 4.63 | 0.74 |
| <i>RPS26A</i> | 4.36 | 0.80 | <i>rplI</i> | 4.68 | 0.73 | <i>rpL20</i> | 4.68 | 0.73 |
| <i>RPS20</i> | 4.39 | 0.84 | <i>rplW</i> | 4.68 | 0.67 | <i>rpL21</i> | 4.68 | 0.74 |
| <i>RPL30</i> | 4.39 | 0.87 | <i>rplA</i> | 4.76 | 0.78 | <i>rpL1</i> | 4.69 | 0.85 |
| <i>RPL24B</i> | 4.42 | 0.77 | <i>rpsA</i> | 4.81 | 0.78 | <i>rpL19</i> | 4.71 | 0.79 |
| <i>RPS24A</i> | 4.42 | 0.83 | <i>rplJ</i> | 4.84 | 0.65 | <i>rpL11</i> | 4.74 | 0.78 |
| <i>RPS26B</i> | 4.42 | 0.71 | <i>rplU</i> | 4.85 | 0.65 | <i>rpL24</i> | 4.77 | 0.69 |
| <i>RPS23B</i> | 4.43 | 0.80 | <i>rplX</i> | 4.85 | 0.59 | <i>rplQ</i> | 4.78 | 0.73 |
| <i>RPL24A</i> | 4.43 | 0.76 | <i>rplK</i> | 4.86 | 0.72 | <i>rpL13</i> | 4.79 | 0.79 |
| <i>RPS10B</i> | 4.43 | 0.80 | <i>rplS</i> | 4.86 | 0.65 | <i>rpS9</i> | 4.80 | 0.77 |
| <i>RPS6B</i> | 4.44 | 0.85 | <i>rpsF</i> | 4.86 | 0.69 | <i>rpS1</i> | 4.81 | 0.83 |
| <i>RPS25A</i> | 4.44 | 0.73 | <i>rpsJ</i> | 4.86 | 0.58 | <i>rps12</i> | 4.82 | 0.72 |
| <i>RPS10A</i> | 4.45 | 0.84 | <i>rpsL</i> | 4.87 | 0.67 | <i>rpL18</i> | 4.83 | 0.69 |
| <i>RPS6A</i> | 4.45 | 0.84 | <i>rplP</i> | 4.87 | 0.62 | <i>rpL3</i> | 4.83 | 0.79 |
| <i>RPL22A</i> | 4.45 | 0.89 | <i>rplC</i> | 4.87 | 0.72 | <i>rpS7</i> | 4.85 | 0.77 |
| <i>RPS9B</i> | 4.46 | 0.84 | <i>rpsC</i> | 4.89 | 0.74 | <i>rpL10</i> | 4.86 | 0.73 |
| <i>RPL8B</i> | 4.47 | 0.85 | <i>rpsB</i> | 4.90 | 0.78 | <i>rpS11</i> | 4.87 | 0.75 |
| <i>RPS14A</i> | 4.48 | 0.80 | <i>rplV</i> | 4.90 | 0.56 | <i>rpS14</i> | 4.89 | 0.68 |
| <i>RPL8A</i> | 4.50 | 0.84 | <i>rplM</i> | 4.90 | 0.68 | <i>rpS5</i> | 4.90 | 0.73 |
| <i>RPL18A</i> | 4.50 | 0.82 | <i>rplQ</i> | 4.92 | 0.56 | <i>rpL16</i> | 4.91 | 0.75 |
| <i>RPL17A</i> | 4.50 | 0.81 | <i>rplR</i> | 4.93 | 0.62 | <i>rpL2</i> | 4.92 | 0.77 |
| <i>RPL13B</i> | 4.51 | 0.75 | <i>rpsE</i> | 4.94 | 0.60 | <i>rpS8</i> | 4.92 | 0.75 |
| <i>RPS17A</i> | 4.51 | 0.81 | <i>rplD</i> | 4.95 | 0.70 | <i>rpL4</i> | 4.93 | 0.73 |
| <i>RPL36A</i> | 4.51 | 0.68 | <i>rplF</i> | 4.96 | 0.62 | <i>rbfA</i> | 4.93 | 0.58 |
| <i>RPL25</i> | 4.52 | 0.76 | <i>rpsN</i> | 4.96 | 0.55 | <i>rpS10</i> | 5.01 | 0.68 |
| <i>RPL19B</i> | 4.52 | 0.72 | <i>rplB</i> | 4.97 | 0.72 | <i>rpS2</i> | 5.02 | 0.78 |

level of confidence were found to be -0.672 , -0.669 and -0.563 for *E. coli*, *S. cerevisiae* and *H. influenzae* respectively. From Figure 1 it can be concluded that the uncertainty values are highly correlated with gene expression levels (CAI). One should keep in mind that the uncertainty values were calculated for all the genomes under study using the simple Shannon's uncertainty formula as described earlier, but to calculate CAI one would require a different reference subset of highly expressed genes from their respective genomes.

It is also interesting to note that the uncertainty value ranges from 3.99 to 5.73 for *E. coli* genes, 4.05 to 5.77 for *S. cerevisiae* genes and 4.10 to 5.57 *H. influenzae* genes. We found no genes with uncertainty value of 5.93 bits per codon, indicating that all genes have some degree of internal constraints in their codon usage. One should not view the low correlation values for *H. influenzae* seriously, since uncertainty values try to capture all the hidden constraints (bias) that nature would face in getting the gene better-adapted to a given genome, while CAI values were calculated from the codon usage of optimal codons of a given gene with respect to a highly expressed gene present in the reference subset for the respective genome.

The low correlation value between the uncertainty and CAI values in the case of *H. influenzae* may be due to the fact that mutational bias is more effective than translational selection in detecting the codon bias in this highly skewed organism, particularly in the presumably highly expressed genes of ribosomal genes.

Genes are sorted according to the uncertainty values. Tables 1 and 2 show the uncertainty and their corresponding CAI values for 50 genes for each organism located at the two extreme ends in their uncertainty scale. From Tables 1 and 2 it is obvious that there is a high degree of correlation between uncertainty values and their corresponding CAI values in *E. coli* and *S. cerevisiae* genes located at the extreme end of the lowest uncertainty values. From Table 2 it was observed that the CAI for *H. influenzae* is comparatively higher than the *E. coli* and *S. cerevisiae* counterparts. The comparatively higher value of CAI of *H. influenzae* may be due to the fact that CAI values for its individual genes are related to its base composition, rather than its expression level. To find out whether genomic GC has any effect on uncertainty values, we have calculated the uncertainty values for the same functional category of genes such as ribosomal proteins of *E. coli*, *S.*

Table 4. Uncertainty values in bits per codon for experimentally known highly expressed genes of *S. cerevisiae*

| Gene | Uncertainty in bits per codon |
|---------------|-------------------------------|
| <i>ENO2</i> | 4.52 |
| <i>TDH3</i> | 4.53 |
| <i>RPS31</i> | 4.53 |
| <i>RPL1A</i> | 4.54 |
| <i>RPL27A</i> | 4.56 |
| <i>TDH2</i> | 4.57 |
| <i>PGK1</i> | 4.60 |
| <i>PDC1</i> | 4.60 |
| <i>RPL2A</i> | 4.60 |
| <i>RPL5</i> | 4.60 |
| <i>FBA1</i> | 4.61 |
| <i>GPM1</i> | 4.61 |
| <i>TEF2</i> | 4.61 |
| <i>TEF1</i> | 4.62 |
| <i>RPL4A</i> | 4.64 |
| <i>ADH1</i> | 4.65 |
| <i>RPL16A</i> | 4.69 |
| <i>RPS18B</i> | 4.70 |
| <i>ADH2</i> | 5.10 |

cerevisiae and *H. influenzae* (Table 3). It was observed that the uncertainty values range between 3.99 and 5.02 bits per codon, emphasizing that majority of the genes of the same functional categories share the same degree of constraint irrespective of their genomic GC composition.

Uncertainty values of experimentally known highly expressed genes of *S. cerevisiae* were calculated to compare them with the other highly expressed genes based on CAI values (Table 4). It is evident from Table 4 that all the experimentally known highly expressed genes fall in the range of 4.52 to 5.09 bits per codon. These values fall in the range of other highly expressed genes, as shown for the ribosomal proteins in Table 3.

From these results we can conclude that Shannon's uncertainty principle can be used as an alternative method for predicting the level of gene expression. In other words, it can be said that Shannon's uncertainty values determine the degree of constraint present in the given gene, which is strongly correlated to the level of gene expression. This method does not require any prior reference subset of highly expressed genes, as required to calculate CAI values. To have an accurate CAI value, one would require knowing all highly expressed genes of the genome, whereas Shannon's uncertainty is immune to this limitation. For the vast amount of sequenced gene data present for different genomes, one can use the uncertainty values to determine the level of gene expression without knowing their proper functional annotation.

- Ohkubo, S., Muto, A., Kawachi, Y., Yamao, F. and Osawa, S., *Mol. Gen. Genet.*, 1987, **210**, 314–322.
- Wright, F. and Bibb, M. J., *Gene*, 1992, **113**, 55–65.
- Gupta, S. K., Bhattacharyya, T. K. and Ghosh, T. C., *Indian J. Biochem. Biophys.*, 2002, **39**, 35–48.
- Ikemura, T., *J. Mol. Biol.*, 1981, **146**, 1–21.
- Ikemura, T., *J. Mol. Biol.*, 1982, **158**, 573–587.
- Bennetzen, J. L. and Hall, B. D., *J. Biol. Chem.*, 1982, **257**, 3026–3031.
- Gupta, S. K. and Ghosh, T. C., *Gene*, 2001, **273**, 63–70.
- Gouy, M. and Gautier, C., *Nucleic Acids Res.*, 1982, **10**, 7055–7074.
- Andersson, S. G. and Sharp, P. M., *Microbiology*, 1996, **142**, 915–925.
- Ghosh, T. C., Gupta, S. K. and Majumdar, S., *Int. J. Parasitol.*, 2000, **30**, 715–722.
- McInerney, J. O., *Proc. Natl. Acad. Sci.*, 1998, **95**, 10698–10703.
- Chiapello, H., Ollivier, E., Landes-Devauchelle, C., Nitschke, P. and Risler, J.-L., *Nucleic Acids Res.*, 1999, **27**, 2848–2851.
- Kerr, A. R. W., Peden, J. F. and Sharp, P. M., *Mol. Microbiol.*, 1997, **25**, 1177–1179.
- Miranda, A. B. de., Valin, F. A., Jabbari, K., Degraeve, W. M. and Bernardi, G., *J. Mol. Evol.*, 2000, **50**, 45–55.
- Lio, P., Ruffo, S. and Buiatti, M., *J. Theor. Biol.*, 1994, **171**, 215–223.
- Gupta, S. K., Majumdar, S., Bhattacharya, T. K. and Ghosh, T. C., *Biochem. Biophys. Res. Commun.*, 2000, **269**, 692–696.
- D'Onofrio, G., Ghosh, T. C. and Bernardi, G., *Gene*, 2002, **300**, 179–187.
- Adzhubei, A. A., Adzhubei, I. A., Krashennnikov, I. A. and Neidle, S., *FEBS Lett.*, 1996, **399**, 78–82.
- Tao, X. and Dafu, D., *FEBS Lett.*, 1998, **434**, 93–96.
- Sharp, P. M. and Li, W.-H., *Nucleic Acids Res.*, 1987, **15**, 1281–1295.
- Schneider, T. D., Stormo, G. D., Gold, L. and Ehrenfeucht, A., *J. Mol. Biol.*, 1986, **188**, 415–431.
- Hamming, R. W., Entropy and Shannon's first theorem. In *Coding and Information Theory*, Prentice-Hall Inc., Englewood Cliffs, New Jersey, 1980, p. 107.
- Shannon, C. E., *Bell Syst. Tech. J.*, 1948, **27**, 279–423, 623–656.
- Velculescu, V. E. *et al.*, *Cell*, 1997, **88**, 243–251.

ACKNOWLEDGEMENTS. We thank the Department of Biotechnology, Government of India for funding the BTIS programme at Bose Institute. We also thank Mr J. N. Mandal for help in preparing the manuscript.

Received 28 August 2003; revised accepted 12 December 2003

1. Wada, K., Aota, S., Tsuchiya, R., Ishibashi, F., Gojobori, T. and Ikemura, T., *Nucleic Acids Res. (Suppl.)*, 1990, **18**, 2367–2411.