

J. Biosci., Vol. 20, Number 5, December 1995, pp 613-627. © Printed in India.

Analysis of CAG/CTG triplet repeats in the human genome: Implication in transcription factor gene regulation

RASHNA BHANDARI and SAMIR K BRAHMACHARI*

Molecular Biophysics Unit, Indian Institute of Science, Bangalore 560 012, India

MS received 15 September 1995; revised 14 November 1995

Abstract. Instability and polymorphism at several CAG/CTG trinucleotide repeat loci have been associated with human genetic disorders. In an attempt to identify novel sites that may be possible loci for expansion of CAG/CTG repeats, we searched all human sequences in the EMBL nucleotide sequence database for (CAG)₅ and (CTG)₅ repeats. We have identified 121 human DNA sequences of known and unknown functions that contain stretches of five or more CAG or CTG repeats. Many repeat stretches were interrupted by variant triplets, a significant number of which differ from the repeat triplet only by a single base, suggesting that these evolved from the parent triplet by point mutations. A large number of human transcription factor genes were found to contain CAG repeats within their coding sequences. Analysis of the EMBL transcription factors database showed that many transcription factor genes of other eukaryotes, including genes involved in *Drosophila* embryo development, possess these repeats. Interestingly, CAG repeats are absent from prokaryotic transcription factors. Different sequence entries for the human TATA box binding protein showed a polymorphism in the length of the CAG repeat in this gene, suggesting that loci other than those already known to be associated with genetic diseases may be possible sites for repeat instability related disorders. On the basis of our findings in this database analysis, we propose a role for CAG repeats as cis-acting regulatory elements involved in fine-tuning gene expression.

Keywords. Triplet repeats; sequence analysis; transcription factor genes; TBP.

1. Introduction

Trinucleotide repeat regions in the genome have recently been shown to be unstable and polymorphic. The expansion of trinucleotide repeats in humans has been associated with several human genetic diseases (Sutherland and Richards 1995). In the Fragile X syndrome, there is expansion of a CGG repeat in the 5'UTR of the FMR-1 gene (Kremer *et al* 1991). Myotonic dystrophy (DM) is associated with the expansion of a CTG repeat in the 5'UTR of the myotonin kinase gene (Mahadevan *et al* 1992; Fu *et al* 1992). Six neurodegenerative disorders—Huntington's disease (HD; Huntington's Disease Collaborative Research Group 1995), spinal and bulbar muscular atrophy (SBMA; La Spada *et al* 1991), spinocerebellar ataxia type 1 (SCA-1; Orr *et al* 1995), dentatorubral pallidoluysian atrophy (DRPLA; Koide *et al* 1994), Haw river syndrome (Burke *et al* 1994), and Machado-Joseph disease (MJD; Kawaguchi *et al* 1994)—are caused by an amplified CAG repeat

Corresponding author (Fax: 91-080-3341683; Email: skb@mbu-iisc.ernet.in).

within the coding sequence of the gene, resulting in an increase in the length of a polyglutamine tract in the corresponding protein. Despite a variation in the location of the repeat within the coding region, all these inherited diseases share the characteristics of increased severity with increasing repeat length. The age of onset of disease decreases in subsequent generations, thus exhibiting genetic anticipation (Richards and Sutherland 1992).

With the discovery that trinucleotide repeat expansions are involved in genetic diseases, attempts have been made to identify other repeat regions that may be possible sites for similar expansions. The screening of cDNA libraries has revealed that CAG and CGG repeats are found in close proximity to several gene loci (Riggins *et al* 1992; Li *et al* 1993). By computer aided sequence analysis of approximately 10 Mb human DNA sequence in GenBank, Han *et al* (1994) have studied the frequency distribution of all the 64 possible trimers. They found that the disease associated trimers, GGC and CAG are over-represented in the genome. They also identified 51 human genes that contain $(GGC)_{\geq 4}$ repeats. A search carried out by Stallings (1994) on approximately 16 Mb GenBank human DNA sequences for all 10 possible microsatellite motifs with a perfect string of eight or more trinucleotides, showed that CAG is the most abundant triplet repeat unit in the human genome.

There are, to date, no composite data available on the location of CAG repeats in the human genome, except for the loci known to be associated with genetic diseases. It is possible that there are other regions in the genome containing long stretches of CAG/CTG repeats that show length polymorphism. The identification and analysis of genes associated with such regions may provide a clue to understanding the functional significance of these repetitive DNA sequences in the genome and perhaps also of the genetic diseases associated with them. In an attempt to identify other sites in the human genome that may be possible loci for expansion of CAG/CTG repeats, we carried out a search on the EMBL nucleotide sequence database. All human sequences in the EMBL database (Release 41, Dec. 1994) were screened for the presence of a stretch of at least 5 tandem repeats of CAG or CTG. We observed that a large number of human transcription factor genes contain $(CAG)_{\geq 5}$ repeats in their coding regions. Analysis of the transcription factors of other organisms listed in the EMBL transcription factors database showed that this repeat is present in many eukaryotic transcription factor genes, but is absent from the genes encoding prokaryotic transcription factors. Based on our analysis, we have examined the potential role of CAG/CTG repeats as *cis*-acting regulatory elements that quantitatively modulate gene expression.

2. Materials and methods

In order to identify human sequences containing CAG or CTG repeats, the FASTA program was used to search all primate sequences in the EMBL nucleotide sequence database (Release 41, Dec. 1994) for the strings, $(CAG)_5$ and $(CTG)_5$. Since the FASTA program searches for a string only once in a sequence, the sequences containing a $(CAG)_5$ or $(CTG)_5$ repeat were put through the GCG Find program (Devereux *et al* 1984) to search for $(CAG)_5$ strings elsewhere in the same sequence.

The annotation of each sequence containing a $(CAG)_{\geq 5}$ or $(CTG)_{\geq 5}$ repeat was then examined to determine whether the repeat occurred in the 5' untranslated

region (5'UTR), coding DNA sequence (CDS), 3'UTR, or in an intron. If a repeat was present in the coding sequence, the amino acid coded for by the triplet repeat was determined by a simple calculation using the start position of the coding sequence given in the feature table of the EMBL sequence entry. The location of the repeat region could not be determined for unannotated sequences and sequences for which the entire CDS was not defined.

The sequence in the vicinity of each stretch of five repeats was then examined. In several sequences we found clusters of CAG/CTG triplets around the tandem (CTG/CTG)_{≥5} repeat. CAG/CTG triplets in this region, that were separated by not more than 12 nucleotides, were included in the repeat stretch shown in table 1. The ends of the repeat were defined by the presence of at least two tandem CAG/CTG triplets that were separated from the nearest CAG/CTG by not more than 12 nucleotides. In most cases the CAG/CTG triplets were interrupted by multiples of three nucleotides. The sequences of variant triplets (T) that interrupted the perfect repeats were noted and the frequency of occurrence of each interrupting triplet was determined.

A large number of human transcription factor genes were found to contain (CAG)_{≥5} stretches. To examine the transcription factors of other organisms, the EMBL transcription factors database was searched for (CAG)₅ stretches using the GCG Find program. All non human transcription factors containing at least five tandem CAG/CTG repeats were examined for repeat and non-repeat triplets in that region, as above.

3. Results

3.1 *Distribution of CAG repeats in human DNA sequences*

We have screened approximately 54 Mb DNA from 52549 primate sequence entries in the EMBL nucleotide sequence database and have identified 121 unique sequences in the human genome that contain (CAG/CTG)_{≥5} repeats. The observed number of repeats is significantly higher than expected [$P \ll 0.01$, based on Poisson distribution, with mean number of (CAG/CTG)₅ repeats per sequence = 10^{-6}]. Table 1 lists the human sequences that contain five or more tandem repeats of CAG/CTG. The subscript denotes the length of CAG(R) and CTG(R') repeats found in these sequences. For those sequence entries in which the CDS is defined, we have determined the location of the repeat with respect to the coding sequence. The repeat CAG codes for Gln, Ser, or Ala in the three reading frames. Similarly, the repeat CTG codes for Leu, Cys, or Ala. In 22 of the 57 sequences in which the reading frame of CAG repeats present within the coding region could be determined, the repeat CAG codes for a stretch of Gln residues. In ten cases, the reading frame is AGC, coding for Ser, and in five, GCA, encoding Ala. The CTG repeat codes for Leu in all cases.

From among the 48 (CAG)_{≥5} repeats in human sequences for which the location with respect to the transcribed region could be determined, it was found that none were located within an intron. However three (CTG)_{≥5} repeats (i.e., CAG on the opposite strand) out of 44 were found within introns. In other words, CAG may be excluded from introns of human genes in a strand specific manner, with no

Table 1. Human sequence entries in the EMBL database that contain (CAG/CTG)_{≥5} repeats.

Entry name	Description	Repeat	Location	Amino acid ^a
HS02013*	Sterol regulatory element binding protein	R ₂ TR ₂ TR ₆ T ₂ RTR ₃	CDS	Ser
HS04840	Onconeural ventral antigen-1 (Nova-1)	R ₅	CDS	Ala
HS07857	18 kDa Alu RNA binding protein	R ₈ T ₃ R ₂ TRT ₃ R ₂ T ₂ R ₂ TR ₂	CDS	Ala
HS9515	T-lymphocyte repeat containing mRNA	R ₅	—	—
HSAF9X*	AF-9	R ₈ TR ₁₀ TR ₁₀ TR ₃ TR ₆	CDS	Ser
HSAFPEBP*	α-fetoprotein enhancer binding protein	R ₇ T ₂ RTR ₂	CDS	Gln
		R ₅ T ₃ R ₄	CDS	Gln
HSALMG1	Calmodulin gene exon 1 (CALM 1)	R ₇	—	—
HSANDREC* ^{b,c}	Androgen receptor	R ₂₀ N ₁₈ R ₆	CDS	Gln
HSASH1A*	Achaete scute homologous protein (ASH 1)	R ₁₄	CDS	Gln
HSAUTANT	Auto antigen	R ₅ T ₂ R ₆ TR ₂	CDS	Ser
HSBNPA	Brain natriuretic peptide	R ₇ TR ₄	5'UTR	—
HSCDP*	CCAAT displacement protein	R ₅	—	—
HSCEI5A	CEI5, 3'flank	R ₆	3'UTR	—
HSCOUPII*	Chick ovalbumin upstream promoter transcription factor II	R ₅	5'UTR	—
HSCYT2A	Cytokeratin 2	R ₇	CDS	Ser
HSDB11*	DB1 (IL-3 promoter binding zinc finger protein)	R ₇	CDS	Gln
HSDNAPLYC	Repeat region DNA (wgl1d10)	R ₈	—	—
HSDRPLA	DRPLA mRNA	R ₁₀	CDS	Gln
HSEGR1* ^c	Early growth response protein 1	R ₆ N ₃₃ R ₆	CDS	Ser
HSEGR2A*	Early growth response protein 2	R ₅	CDS	Ala
HSENKB4	Preproenkephalin B	R ₂ TR ₇	3'UTR	—
HSEPKER	Type I epidermal keratin	R ₆	CDS	Ser
HSERF2	ERF-2	R ₇	CDS	Gln
HSERK1	ERK1 (protein serine threonine kinase)	R ₆	5'UTR	—
HSFRI	Chromosome 13q12 DNA fragment	R ₅	—	—
HSHBRM*	hbrm (homologue of <i>Drosophila</i> brm)	R ₃ TR ₁₃ TR ₂ T ₄ RTRT ₂ R ₂	CDS	Gln
HSHCCA	Splicing factor (CC1.3)	R ₇	5'UTR	—
HSHSNFA*	hSNF2a (estrogen receptor)	R ₃ TR ₁₃ TR ₂ TR ₃	—	—
HSI9R	Interleukin 9 receptor	R ₂ TR ₈	CDS	Ser
HSIFNRTF*	IFN responsive transcription factor	R ₆	CDS	Ser
HSJUNA*	c-jun proto-oncogene	R ₅	CDS	Gln
HSKER16A1	Keratin type 16 gene, exon 1	R ₅	—	—
HSKERP1	Keratin psuedogene, exon 1	R ₅	—	—
HSMEF2*	Myocyte-specific enhancer factor 2	R ₁₁	CDS	Gln
HSMMDA	Human DNA from cosmid MMDA from chromosome 19q13	R ₅	—	—
HSMRNAC	Human brain cDNA (clone CTG-A4)	R ₄ TR ₉ TRTR ₆ T ₂ R ₂	—	—
HSMRNAD1	Human brain cDNA (clone CTG-B10)	R ₅ TR ₉	—	—

Table 1. (Contd)

Entry name	Description	Repeat	Location	Amino acid ^a
HSMYC1*	c-myc oncogene	R ₅	CDS	Gln
HSNATPEP	Natriuretic peptide precursor	R ₇ TR ₄	5'UTR	—
HSNFYA*	CAAT box DNA binding protein subunit A	R ₅	CDS	Gln
HSNOCTPOU*	N-Oct3, N-Oct5a and N-Oct5b	R ₃ TR ₃ TRTR ₃ TR ₂ TR ₆	CDS	Gln
HSNOTIC	Human chromosome 13q12-specific mRNA	R ₅	—	—
HSNUCLEOB*	Nucleobindin	R ₆	CDS	Gln
HSOFRIA	ORF of a gene expressed in human hypothalamus	R ₁₄	CDS	Ser
HSORFG	Human mRNA for randomly sequenced product	R ₂ TR ₅ TRN ₈ RN ₁₂ R ₂	CDS	Ala
HSPDEE	Phosphodiesterase	R ₅	5'UTR	—
HSPIM1A	pim-1 proto-oncogene	R ₈	5'UTR	—
HSREPRA	Simple repeat region (wg1a9)	R ₃ T ₃ R ₁₂	—	—
HSRNAF	Brain cDNA (clone CTG-B37)	R ₁₅	—	—
HSRNAG	Brain cDNA (clone CTG-B43A)	R ₃ TR ₅ TR ₂ TR ₅ TR ₁₀	—	—
HSRSRFC49*	SRF-related proteins RSRFC4 and RSRFC9	R ₉	CDS	Gln
HSRWG1A10	Simple DNA sequence region (clone wg1a10)	R ₅	—	—
HSRWG1A4	Simple DNA sequence region (clone wg1a4)	R ₁₀	—	—
HSRWG1A5	Simple DNA sequence region (clone wg1a5)	R ₉	—	—
HSRWG1F9	Simple DNA sequence region (clone wg1f9)	R ₉	—	—
HSSATB1A*	MAR/SAR DNA binding protein	R ₂ N ₂ RTRTR ₃ TRTR ₇	CDS	Gln
HSSATM01	DNA sequence 3'-flanking mini-satellite pMS1	R ₆	—	—
HSSCA1	Ataxin (SCA-1)	R ₁₂ TRTR ₁₅ N ₅ R ₂	CDS	Gln
HSSRP14A	Signal recognition particle subunit 14	R ₈ T ₃ R ₂ TR ₂ T ₂ R ₂ T ₂ R ₂ TR ₂	CDS	Ala
HSTFEB*	Transcription factor TFEB	R ₆ TR ₃	CDS	Gln
HSTFIID* [†]	TATA binding protein (TFIID) (HeLa cells)	R ₃ T ₃ R ₉ TRTR ₁₈ TRNRN ₃ RN ₃ RN ₅ R ₂ N ₉ R ₂	CDS	Gln
HSTFIIDAA* [†]	Transcription factor TFIID (human hypothalamus)	R ₃ T ₃ R ₁₁ TRTR ₁₆ TRNR N ₃ RN ₃ RN ₅ R ₂ N ₉ R ₂	CDS	Gln
HSTFIIDX* [†]	Transcription factor TFIID (Namalwa cell line)	R ₃ T ₃ R ₉ TRTR ₁₄ TRNRN ₃ RN ₃ RN ₅ R ₂ N ₉ R ₂	CDS	Gln
HSTGFBC	Transforming growth factor-β	R ₂ T ₂ R ₆	CDS	Gln
HSTRINUC	Huntington's disease gene (HD)	R ₂₁	CDS	Gln
HSTUPLE1*	TUP1-like enhancer of split gene 1	R ₅	CDS	Gln
HSWGIF4RP	Repeat region DNA (clone wg1f4)	R ₂ N ₄ R ₁₀	—	—

Table 1. (Contd)

Entry name	Description	Repeat	Location	Amino acid ^a
S62539	Insulin receptor substrate-1	R ₈ R ₇ R ₆	5'UTR, CDS CDS	— Ser Gln
A08001	Repeat region	R ₆ '	—	—
A18521	Mullerian inhibiting substance	R ₅ '	CDS	Leu
HS03398	Receptor 4-1BB ligand	R ₆ '	CDS	Leu
HS04806	FLT3/FLK2 receptor tyrosine kinase ligand	R ₅ '	CDS	Leu
HS09510	Glycyl tRNA synthetase	R ₅ '	5'UTR	—
HS10886	Density enhanced phosphatase-1	R ₆ '	CDS	Leu
HS8955	β-spectrin (HSPTB1)	R ₅ '	Intron	—
HSA3NARSP	Nicotine acetylcholine receptor α-3 subunit	R ₆ 'TR ₂ '	CDS	Leu
HSAK1	Cytosolic adenylate kinase (AK1)	R ₇ '	Intron	—
HSALPP	Placental alkaline phosphatase	R ₇ '	CDS	Leu
HSAMYB15	Amyloid-β protein (APP)	R ₆ '	Intron	—
HSAP2*	Transcription factor AP2	R ₃ 'TR ₅ '	3'UTR	—
HSAPB01	Apolipoprotein B-100	R ₂ 'TR ₂ 'TR'T ₂ R ₆ '	CDS	Leu
HSAREB6*	Transcription factor AREB6	R ₇ '	3'UTR	—
HSARSBX	Aryl sulphatase B	R ₅ '	CDS	Leu
HSBPGM1	2,3-bisphosphoglycerate mutase	R ₈ '	5'UTR	—
HSCIINH2	C1 inhibitor	R ₂ 'T ₂ R ₅ '	CDS	Leu
HSCA1V	Collagen α-1 (V) chain	R ₂ 'T ₂ R ₆ '	CDS	Leu
HSCN2	Skin collagenase	R ₅ '	CDS	Leu
HSCP2101	Mutant CYP 21 B (steroid 21 mono-oxygenase)	R ₅ 'TR ₂ 'T ₄ R ₂ '	CDS	Leu
HSD9S118	DNA segment containing GT repeat (D9S118)	R ₈ '	—	—
HSDIPEPA	Dipeptidyl aminopeptidase like protein	R ₅ '	5'UTR	—
HSEFL2	EHK-1 receptor tyrosine kinase ligand (EFL-2)	R ₆ '	CDS	Leu
HSENA78	ENA-78 (neutrophil activating peptide)	R ₅ '	CDS	Leu
HSEPOHYDD	Epoxide hydrolase	R ₆ '	3'flank	—
HSERK1	ERK1 (protein serine threonine kinase)	R ₅ '	3'UTR	—
HSFLA1A	Leucocyte associated molecule-1 α-subunit	R ₅ '	CDS	Leu
HSGAS	Growth-arrest-specific protein	R ₆ '	CDS	Leu
HSGFAP	Glial fibrillary acidic protein	R ₆ '	5'UTR	—
HSGLUT5	Glucose transport like 5 (GLUT5)	R ₅ 'N ₁₁ R' ₉ R ₂ '	CDS	Leu
HSHGFP	Hepatocyte growth factor activator precursor	R ₆ '	CDS	Leu
HSHNF4	Hepatocyte nuclear factor 4	R ₅ '	CDS	Leu
HSIIP	γ-interferon-inducible protein (IP3)	R ₂ 'TR'T ₂ R ₅ '	CDS	Leu

Table 1. (Contd)

Entry name	Description	Repeat	Location	Amino acid ^a
HSKIAAC	ORF	R ₅	CDS	Leu
HSLHHCGR	Leutenising hormone-choriogonadotropin receptor	R ₂ TR ₅	CDS	Leu
HSLPAP	Lymphocyte phosphatase associated protein	R ₅ T ₂ R ₂	CDS	Leu
HSM DAPKB [†]	Myotonic dystrophy associated protein kinase	R ₁₁	3'UTR	—
HSMYOPK [†]	Myotonin protein kinase associated triplet repeat region	R ₅	3'UTR	—
HSNNAR	α-3 neuronal nicotinic acetyl-choline receptor subunit	R ₆ TR ₂	CDS	Leu
HSNUCPHOF	Processed pseudogene of nucleophosmin/B23	R ₇	—	—
HSORFM	Randomly sequenced product	R ₆ T ₃ R ₂	—	—
HSPA1V	pro-α-1 (V) collagen	R ₂ T ₂ R ₆	CDS	Leu
HSPPR	Protective protein	R ₈	CDS	Leu
HSPSPBA	Pulmonary surfactant associated protein PSP-B	R ₂ T ₂ R ₅ N ₈ R ₁₁ R ₂	CDS	Leu
HSREPRC	Simple repeat region (wg1d6)	R ₆	—	—
HSRNAE	Human brain cDNA (clone CTG-B33)	R ₁₃	—	—
HSRWG1E10	Simple DNA sequence region (clone wg1e10)	R ₈	—	—
HSRWG2A11	Simple DNA sequence region clone (wg2a11)	R ₇	—	—
HSRWG2B8	Simple DNA sequence region clone (wg2b8)	R ₂ T ₃ R ₆ N ₄ R ₂	—	—
HSRWG2D12	Simple DNA sequence region clone (wg2d12)	R ₂ T ₂ R ₄ R ₄ TR ₂ R ₂ TR ₂ T ₂ R ₈ T ₄ R ₄	—	—
HSSIRPOCK	Simple repeat polymorphism	R ₇	—	—
HSSUBTDA	Subtelomeric DNA sequence	R ₆ TR ₂ R ₂	—	—
HSTGFBI	Transforming growth factor-β precursor	R ₅ T ₂ R ₂	CDS	Leu
HSUVOECAD	Uvomorulin (E-cadherin)	R ₅	CDS	Leu
HSVLDLR01	Very low density lipoprotein receptor	R ₅	5'flank	—
S57777	N-acetyl galactosamine-4-sulphatase	R ₅	CDS	Leu
S57793	Lutenising hormone receptor	R ₂ TR ₅	CDS	?

R, CAG triplet; R', CTG triplet; T, any triplet; N, single nucleotide.

^aThe amino acid coded for by the (CAG/CTG)₅ repeat is given. The reading frame of neighbouring CAG/CTG triplets may not be the same.

^bDifferent sequence entries for androgen receptor have number of repeats ranging from 15 to 20.

^cTwo repeats found within the same sequence are close enough to be located within the same nucleosome.

*Human transcription factor genes.

[†]Different entries for the same sequence, which differ in the number of triplet repeats.

?, Author given protein sequence is in conflict with the conceptual translation.

CAG repeats being seen in the sense strand. A similar result was obtained by Stallings (1994) from a search on the GenBank DNA database, wherein it was found that none of the 15 (CAG)_{≥8} repeats found in mouse, human or rat sequences was located within intronic sequences. It has been suggested that the exclusion of CAG from introns is possibly because of its similarity to the 3' consensus splice site, CAGG (Stallings 1994).

In 43 sequences we found a stretch of interrupted CAG or CTG repeats. Long stretches of interrupted repeats were seen in the Alu RNA binding protein, the human homologue of *Drosophila* brm, the signal recognition particle subunit 14, the transcription factor N-Oct, the human TATA-binding protein, and in apolipoprotein B-100. A striking feature was that interruptions by multiples of three nucleotides were significantly higher than expected ($\chi^2 = 113.9$ for interruptions of CAG and $\chi^2 = 46.4$ for CTG; $P \ll 0.001$). In most cases, the interruptions were by triplets, such that the reading frame of successive perfect repeats was maintained. Non-triplet interruptions were seen only in eight sequences. The number of intervening nucleotides (N) and triplets (T) within each repeat sequence is shown in table 1.

Table 2 lists the interrupting triplets (T) that were in frame with the (CAG/CTG)_{≥5} repeats, and the number of times each interrupting triplet occurred. By simple statistical analysis it can be seen that if the interrupting triplets had arisen randomly, then only 9/63 (i.e., 14%) of them would differ from CAG or CTG by one base. However, we see that 88% of the interrupting triplets differ from the flanking CAG

Table 2. Frequency distribution of triplets interrupting (CAG/CTG) repeats.

Parent triplet	Type of mutation	Interrupting triplet	Base change	Occurrence	Frequency (%)
CAG	Point mutation (transition)	TAG	C → T	8	9.1
		CGG	A → G	10	11.4
		CAA	G → A	41	46.6
	Point mutation (transversion)	AAG	C → A	2	2.3
		GAG	C → G	3	3.4
		CCG	A → C	5	5.7
		CTG	A → T	2	2.3
		CAC	G → C	3	3.4
		CAT	G → T	3	3.4
	Two or three base changes			11	12.4
	CTG	Point mutation (transition)	TTG	C → T	12
CCG			T → C	5	8.8
CTA			G → A	4	7
Point mutation (transversion)		ATG	C → A	1	1.8
		GTG	C → G	1	1.8
		CAG	T → A	1	1.8
		CGG	T → G	—	—
		CTC	G → C	5	8.8
CTT		G → T	—	—	
Two or three base changes				28	49

triplets by one base change, and in the case of CTG repeats, 51 % of the variant triplets differ by a single base. Since these values are significantly higher than expected, the null hypothesis of "random origin" can be rejected. (Expected frequencies are computed using binomial distribution, $\chi^2 = 384.2$ for CAG interruptions and $\chi^2 = 62.8$ for CTG interrupting triplets; $P \ll 0.001$). The indication is that the interrupting triplets are very likely to have arisen by point mutations in an ancestral CAG or CTG repeat sequence. We also observe that single base transitions, i.e., $A \leftrightarrow G$ or $C \leftrightarrow T$, are more frequent than transversions. Transitions were approximately three-fold more common than transversions in the case of CAG and approximately two-fold for CTG. It has been seen that transitions are two times more likely than transversional mutations (Kimura 1983). The high frequency of CAG to CAA mutations may be due to the synonymous nature of this substitution.

3.2 Several transcription factor genes contain CAG repeats

A large number of human transcription factors were found to contain CAG repeats in their coding sequences. To look for transcription factors in other organisms that contain CAG repeats, sequences listed in the EMBL transcription factors database were searched for strings of $(CAG)_5$. Out of 563 transcription factor gene entries from 86 different species analysed, 48 transcription factors from 12 species were found to contain $(CAG)_{\geq 5}$ repeats (tables 1 and 3). A number of *Drosophila* genes involved in embryo development, such as *dorsal*, *knirps*, *deformed* and *hunchback*, possess long CAG stretches in their coding regions. In most cases, CAG repeats are present in the sense strand of transcription factors. Only in human transcription factors AP2 and AREB6 are CTG stretches found in the coding strand. Most transcription factors monitored contain CAG repeats in the reading frame that codes for glutamine, but some also contain serine and alanine stretches coded for by this repeat. An interesting observation was that although CAG repeats were present in the transcription factors of lower and higher eukaryotes, none of the prokaryotic transcription factor genes contained CAG repeats.

3.3 Length polymorphism in human TATA binding protein

Human TATA binding protein has a stretch of CAG repeats similar in length to loci such DRPLA, MJD, and SCA-1 that show dynamic mutations leading to repeat expansion. It is therefore quite possible that the repeat stretch at the TFIID locus may similarly undergo dynamic mutation. Figure 1 shows the alignment of the repeat region of the MJD 1a and SCA-1 loci with the different sequence entries for human TBP that are derived from different sources, viz., HeLa cells, Namalwa cells, and human hypothalamus. It can be seen that hTBP contains three CAG repeats that are interrupted by the variant triplet CAA, which also codes for glutamine. Similar variations on the CAG motif are also seen in the MJD 1a locus. The 5' end of the CAG repeat in this locus contains two variant sequences, CAA and AGG, at three positions. The positions at which the variant CAA triplets occur is conserved in all the sequence entries of hTBP, except at two positions in the TBP cloned from human hypothalamus.

It is interesting to note that there is a polymorphism in the length of the repeats in TBP derived from different sources. The number of repeats in the 3' CAG

Table 3. Non-human transcription factor genes containing (CAG/CTG)_{≥5} repeats.

Entry name	Description	Species	Repeat	Location	Amino acid
BTTRFSQA	Transcription factor	<i>Bos taurus</i> (cattle)	R ₅	—	—
DM74E	74E	<i>D. melanogaster</i>	R ₅ TR ₂ TRT ₂ R ₂ R ₂ N ₁₀ R ₅	CDS CDS	Gln Ser
DMDFD	Deformed (Dfd)	<i>D. melanogaster</i>	R ₂ T ₂ R ₅ TR ₄ TR ₂ TR ₂	CDS	Gln
DMDORSAL	Dorsal	<i>D. melanogaster</i>	R ₂ T ₃ R ₄ TR ₄ TRTRT RTR ₂ T ₃ R ₉ R ₄ TR ₇	CDS CDS	Gln Gln
DMHGB	Hunchback	<i>D. melanogaster</i>	R ₅ TR ₂	CDS	Gln
DMKNIRPS	Knirps	<i>D. melanogaster</i>	R ₇		
DMTATABF	TFIID (TATA binding factor)	<i>D. melanogaster</i>	R ₅ TR ₂	CDS	Gln
DVHB	Hunchback (hb)	<i>D. virilis</i>	R ₂ TR ₇	CDS	Gln
FCFTT	Feline proviral (FTT) v-myc	<i>Felis catus</i> (cat)	R ₅	CDS	Gln
GDCMD1	CMD1 (mouse myoD1 homologue)	<i>G. domesticus</i> (chicken)	R ₆	CDS	Ser
MAHCRMVC	Fused her and myc genes	<i>Marmota monax</i> (woodchuck)	R ₅	CDS	Gln
MMCJUN	c-jun oncogene	<i>Mus musculus</i> (mouse)	R ₅	CDS	Gln
MMGRLA	Glucocorticoid receptor form A	<i>Mus musculus</i>	R ₈	CDS	Gln
PTCMYC	c-myc proto-oncogene	<i>Pan troglodytes</i> (chimpanzee)	R ₅	CDS	Gln
REFCVMYC	v-myc gene	<i>Feline leukemia virus</i>	R ₅	CDS	Gln
RNANDREC	Androgen receptor	<i>Rattus norvegicus</i> (rat)	R ₄ TRTR ₁₁ TR ₃	CDS	Gln
RNCBFB	CCAAT binding transcription factor-B subunit	<i>R. norvegicus</i>	R ₅	CDS	Gln
RNNGFIA	Nerve growth factor-induced (NGFI) gene	<i>R. norvegicus</i>	R ₉	CDS	Ser
RNRJG9	RJG-9 gene for c-jun	<i>R. norvegicus</i>	R ₅	CDS	Gln
SCADR6	ADR6	<i>S. cerevisiae</i>	R ₆	CDS	Gln
SCHAP2	HAP2	<i>S. cerevisiae</i>	R ₇	CDS	Gln
SCINO4X	INO4	<i>S. cerevisiae</i>	R ₅	CDS	Ser

stretch is 18 in TBP from HeLa cells and 14 in TBP from Namalwa cells (figure 1). The polymorphism could reflect an instability at this locus in the general population. Analysis of TBP from 48 unrelated individuals has shown that the PCR products encompassing the CAG repeat region of this gene are polymorphic in length (Polymeropoulos *et al* 1991). Perhaps a narrow range of length variation is allowed in the TBP gene, but expansion beyond a certain number could be deleterious.

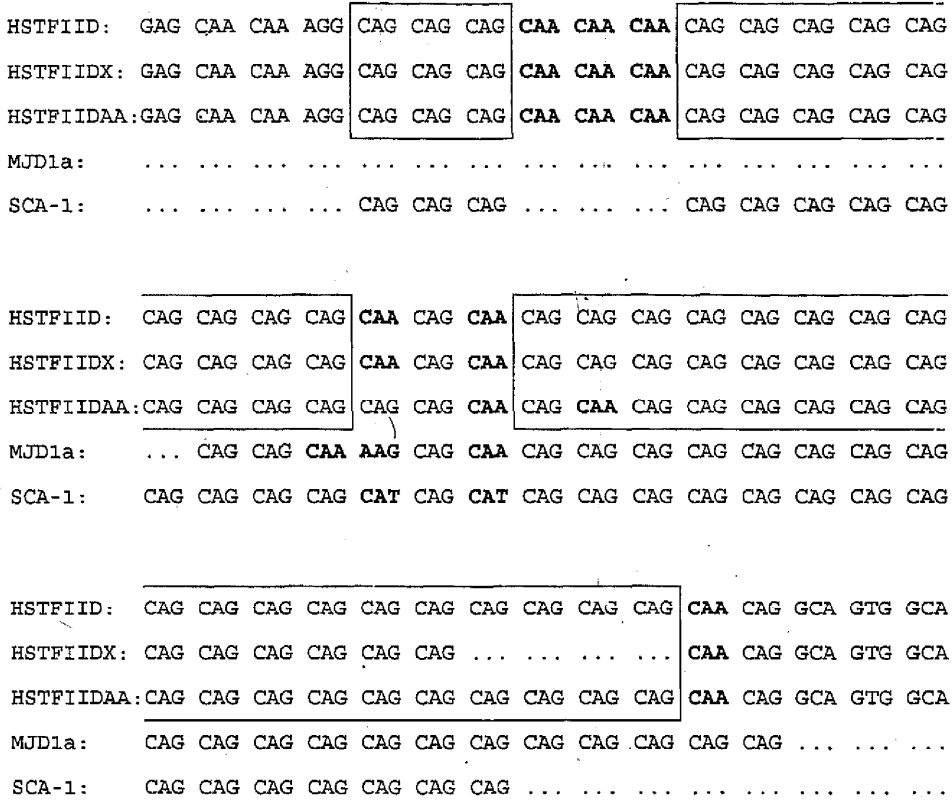


Figure 1. The CAG repeat region of hTBP derived from HeLa cells (HSTFIID), Namalwa cells (HSTFIIDX), and human hypothalamus (HSTFIIDAA), is compared with the CAG repeat of the MJD1a and SCA-1 loci. Gaps have been introduced to align the sequences. The three CAG repeat stretches in hTBP are boxed.

4. Discussion

We have identified 121 human DNA sequences that contain a stretch of five or more CAG or CTG repeats. This list encompasses a variety of different proteins from different tissues, with diverse functions. Apart from the five loci known to be associated with triplet repeat expansion related diseases, i.e., androgen receptor, DRPLA, SCA-1, HD and myotonin kinase, we found 15 loci that have 10 or more tandem CAG/CTG repeats. Of these, 7 are random DNA sequences identified by screening of cDNA libraries (Li *et al* 1995;-Armour *et al* 1994; Phillips *et al* 1995). Other than these, several genes of known function, namely, AF-9, achaete-scute homologous protein, human estrogen receptor, myocyte- specific enhancer factor, human homologue of *Drosophila* brm, and human TATA-box binding protein, contain stretches of (CAG/CTG)_{≥10}. There are also some sequences that contain long stretches of interrupted CAG repeats, which are usually in frame; that is, the repeats are interrupted, in most cases, by triplets.

Earlier studies have suggested that interruption of a repeat stretch serves to stabilize the repeat region. It has been shown that for the SCA-1 locus, 98% of normal alleles show interruption of the CAG repeat by CAT triplets, while disease alleles consist of pure CAG repeats (Chung *et al* 1993). Also, in the case of the FRAXA locus, the CGG repeat is interrupted, at least in normal alleles, by one or two AGG triplets. The DNA sequence in the region containing AGG has been shown to be relatively invariant. Most of the variation at the FRAXA locus occurs at the 3' end, where the longest tracts of uninterrupted CGG are found (Kunst and Warren 1994). In the case of MJD, the CAG repeats contain two variations on the CAG motif at three positions, towards the 5' end of the repeat. CAG expansions in patients are seen to occur on the 3' side of the variant sequences (Kawaguchi *et al* 1994). From figure 1 it can be seen that human TBP has three CAG repeat stretches that are bordered by CAA triplets. It is possible that the presence of CAA triplets prevents expansion of the (CAG)₃ repeat located at the 5' end and the (CAG)₉ stretch in the centre of the repeat region. However, a polymorphism is evident in the long stretch of CAG triplets at the 3' end. The absence of intervening CAA triplets within this stretch may allow variation in repeat length by expansion or contraction during DNA replication. From these observations one can speculate that intervening triplets may serve to stabilize repeat region DNA and preclude expansion of the microsatellite.

Our observation is in agreement with the strand slippage model for expansion of repetitive DNA sequences (Streisinger *et al* 1966). It has been suggested that the loss and gain of repeat units occurs because of errors resulting from strand slippage during DNA replication, which remain uncorrected as a result of defective post replication heteroduplex repair (Strand *et al* 1993). In this picture, if a stretch of CAG repeats is interrupted by other triplets, slippage and looping out of this region would not be allowed, since this would result in extensive mismatch in the loop region and instability of the loop (Kang *et al* 1995). Therefore, expansion would be allowed only at long uninterrupted CAG repeats. It is also observed that the intervening triplets differ from the parent triplet repeat usually by a single base change. A simple explanation for this would be that these interrupting triplets arose during the course of evolution by mutation of the ancestral CAG/CTG triplet, perhaps during incorrect repair following strand slippage, and that these mutants were selected for because they stabilized the length of the repeat region.

Many genes encoding transcription factors have been found to contain long CAG/CTG repeats. It is known that the transcriptional activation domains of many transcription factors contain polyglutamine tracts, and that the length of these tracts may modulate transcription factor activity (Gerber *et al* 1994). It is postulated that these amino acid tracts serve as interfaces for intermolecular protein-protein interactions during transcription activation (Hoffmann *et al* 1990; Stott *et al* 1995). However, our finding that the CAG repeat motif occurs in all three reading frames in transcription factors, implies that stretches of polyserine, and polyalanine are also encoded by this repeat. Therefore in addition to the involvement of CAG encoded polyamino acid tracts in protein-protein interactions, the CAG repeat motif may also play a regulatory role at the DNA level.

It has been hypothesised that tandem repeats may modulate gene expression by variation in their copy numbers (Trifonov 1989; Tripathi and Brahmachari 1991). We have observed that CAG repeats are found only in the transcription factors of

eukaryotes, whose DNA is organized into nucleosomes, and are not found in prokaryotic transcription factors. This suggests that CAG repeats may play a role in histone-DNA interactions. Wang *et al* (1994) have shown that long CTG repeats may act as strong nucleosome positioning signals. It has been shown that expansion of a CTG repeat alters the DNAase I hypersensitivity of the repeat region, perhaps due to a change in nucleosome positioning (Otten and Tapscott 1995). These findings suggest that CAG repeats may serve as nucleosome phasing signals, thereby determining promoter accessibility, and thus regulating transcription of the gene. It is possible that length polymorphism of the repeat is able to affect gene expression by changes in nucleosome affinity. Certain genes such as transcription factors require finely regulated levels of expression. Down regulation of transcription may be achieved by increased binding to histones due to the presence of CAG repeats in these genes, and this may be reversed by histone phosphorylation.

Neurodegenerative disorders that are caused by expansion of CAG repeats within the coding region involve a loss of function, which in some cases has been correlated with decreased mRNA levels (Fu *et al* 1993). Experimental work carried out in our laboratory has shown that the presence of CTG repeats within a reporter gene results in a decrease in the level of gene expression in *Escherichia coli* and yeast, along with a decrease in mRNA level (Brahmachari *et al* 1995). Duplex DNA fragments containing the CTG repeat were seen to exhibit enhanced electrophoretic mobility with increase in repeat length, suggesting that adoption of a compact structure could be responsible for decreased gene expression *in vivo*. Spectroscopic and gel electrophoretic studies carried out in our laboratory on single stranded CAG and CTG repeat sequences show that these oligonucleotides adopt a compact structure. Mitas *et al* (1995) have shown that a single stranded oligonucleotide containing 15 CTG repeats forms a hairpin structure. In addition, a number of RNA binding proteins, some of which recognize hairpin structures, have been identified in eukaryotic systems (Chen and Frankel 1994). Richards *et al* (1993) have shown the existence of a protein that specifically binds CGG repeats. It is quite possible that there may be a similar CAG repeat binding protein. Formation of secondary structure at the DNA and/or mRNA level, and interaction with specific binding proteins may regulate the efficiency of transcription and/or translation of genes containing CAG/CTG repeats.

In the light of our findings in this database search, and experimental evidence, we hypothesise that apart from the involvement of CAG encoded poly amino acid tracts in protein function, the CAG repeat motifs also play a critical role as *cis*-acting regulatory elements at the DNA and the mRNA level. Increased histone binding of CAG containing regions may serve to autoregulate the expression of genes in which they are present. We find that long CAG/CTG repeats occur at many loci, such as the estrogen receptor, myocyte specific enhancer factor, AF-9, TATA binding protein, and in many sequences of unknown function. In addition, polymorphism is observed at the TBP locus. This suggests that in addition to the loci known to be associated with disease, CAG repeats may show polymorphism at many other loci. Slight variations in repeat length at these regions may lead to quantitative differences in the level of expression of these genes, which in turn may be manifested as subtle or overt differences in phenotype. It is possible that the loci at which we have now shown the existence of triplet repeats may be associated with repeat length expansion related diseases that are as yet unidentified.

Acknowledgements

The authors thank Ms. Shyamla at Distributive Information Centre (IISc) and Harish Rao for their assistance in database analysis. We thank Dr N V Joshi, Centre for Ecological Sciences (IISc), for his invaluable help in statistical analysis of the data. SKB wishes to thank Department of Biotechnology, New Delhi, for financial support. We acknowledge the help rendered by T R Raghunand and S Unniraman in data tabulation and compilation. Our thanks are due to S Raghavan and U S Shaligram for their helpful suggestions and critical comments.

References

- Armour J A L, Neumann R, Gobert S and Jeffreys A J 1994 Isolation of human simple repeat loci by hybridisation selection; *Hum. Mol. Genet.* **3** 599-605
- Brahmachari S K, Meera G, Sarkar P S, Balagurumoorthy P, Tripathi J, Raghavan S, Shaligram U S and Pataskar S S 1995 Functional significance of simple repetitive sequences in the genome; *Electrophoresis* **16** 1705-1714
- Burke J R, Kingfield M S, Lewis K E, Roses A D, Lee J E, Huette C, Pericak-Vance M A and Vance J M 1994 The Haw River Syndrome: Dentatorubropallidolusian atrophy (DRPLA) in an African-American family; *Nature Genet.* **7** 521-524
- Chen L and Frankel A D 1994 An RNA-binding peptide from bovine immunodeficiency virus Tat protein recognises an unusual RNA structure; *Biochemistry* **33** 2708-2715
- Chung M Y, Ranum L P, Duvick L A, Servadio A, Zoghbi H Y and Orr H T 1993 Evidence for a mechanism predisposing to intergenerational CAG repeat instability in spinocerebellar ataxia type I; *Nature Genet.* **5** 254-258
- Devereux J, Haerberli P and Smithies O 1984 A comprehensive set of sequence analysis algorithms for the VAX; *Nucleic Acids Res.* **12** 387-395
- Fu Y H, Friedman D L, Richards S, Pearlman J A, Gibbs R A, Pizzuti A, Ashizawa T, Perryman M B, Scarlato G, Fenwick R G and Caskey C T 1993 Decreased expression of myotonin-protein kinase messenger RNA and protein in adult form of myotonic dystrophy; *Science* **260** 235-238
- Fu Y H, Pizzuti A, Fenwick R G Jr, King J, Rajnarayan S, Dunne P W, Dubell J, Nasser G A, Ashizawa T, Jong P D, Wieringa B, Korneluk R, Perryman M B, Epstein H F and Caskey C T 1992 An unstable triplet repeat in a gene related to myotonic muscular dystrophy; *Science* **255** 1256-1258
- Gerber H P, Seipel K, Georgiev O, Hofferer M, Hug M, Rusconi S and Schaffner W 1994 Transcriptional activation modulated by homopolymeric glutamine and proline stretches; *Science* **263** 808-811
- Han J, Hsu C, Zhu Z, Longshore J W and Finley W H 1994 Over-representation of the disease associated (CAG) and (CGG) repeats in the human genome; *Nucleic Acids Res.* **22** 1735-1740
- Hoffmann A, Sinn E, Yamamoto T, Wang J, Roy A, Horikoshi M and Roeder R G 1990 Highly conserved core domain and unique N terminus with presumptive regulatory motifs in human TATA factor (TFIID); *Nature (London)* **346** 387-390
- Huntington's Disease Collaborative Research Group 1993 A novel gene containing a trinucleotide repeat that is expanded and unstable on Huntington's disease chromosomes; *Cell* **72** 971-983
- Kang S, Jaworski A, Ohshima K and Wells R D 1995 Expansion and deletion of CTG repeats from human disease genes are determined by the direction of replication in *E. coli*; *Nature Genet.* **10** 213-218
- Kawaguchi Y, Okamoto T, Taniwaki M, Aizawa M, Inoue M, Katayama S, Kawakami H, Nakamura S, Nishimura M, Akiguchi I, Kimura J, Narumiya S and Kakizuka A 1994 CAG expansions in a novel gene for Machado-Joseph disease at chromosome 14q32.1; *Nature Genet.* **8** 221-228
- Kimura M 1983 *The neutral theory of molecular evolution* (Cambridge: Cambridge Univ. Press)
- Koide R, Ikeuchi T, Onodera O, Tanaka A H, Igarashi S, Endo K, Takahashi H, Kondo R, Ashikawa A, Hayashi T, Saito M, Tomoda A, Miike T, Naito H, Ikuta F and Tsuji S 1994 Unstable expansion of CAG repeat in hereditary dentatorubral-pallidolusian atrophy (DRPLA); *Nature Genet.* **6** 9-12

- Kremer E J, Pritchard M, Lynch M, Yu S, Holman K, Baker E, Warren S T, Schlessinger D, Sutherland G R and Richards R I 1991 Mapping of DNA instability at the fragile X to a trinucleotide repeat sequence p(CCG)_n; *Science* **252** 1711-1714
- Kunst C B and Warren S T 1994 Cryptic and polar variation of the fragile X repeat could result in predisposing normal alleles; *Cell* **77** 853-861
- La Spada A R, Wilson E M, Lubahn D B, Harding A E and Fishbeck K H 1991 Androgen receptor gene mutations in X-linked spinal and bulbar muscular atrophy; *Nature (London)* **352** 77-79
- Li S H, McInnes M G, Margolis R L, Antonarakis S E and Ross C A 1993 Novel triplet containing genes in human brain: Cloning, expression, and length polymorphism; *Genomics* **16** 572-579
- Mahadevan M, Tsilfidis C, Sabourin L, Shutler G, Amemeiya C, Jansen G, Neville C, Narang M, Barcelo J, O'Hoy K, Leblund S, Earle-Macdonald J, deJong P J, Wieringa B and Komeluk R G 1992 Myotonic dystrophy mutation: an unstable CTG repeat in the 3' untranslated region of the gene; *Science* **255** 1253-1255
- Mitas M, Yu A, Dill J, Kamp T J, Chambers E J and Haworth I S 1995 Hairpin properties of single-stranded DNA containing a GC-rich triplet repeat: (CTG)₁₅; *Nucleic Acids Res.* **23** 1050-1059
- Orr H T, Chung M, Banfi S, Kwiatkowski T J, Servadio A, Beaudet A L, McCall A E, Duvick L A, Ranum L and Zoghbi H Y 1993 Expansion of an unstable trinucleotide CAG repeat in spinocerebellar ataxia type 1; *Nature Genet.* **4** 221-226
- Otten A D and Tapscott S J 1995 Triplet repeat expansion in myotonic dystrophy alters the adjacent chromatin structure; *Proc. Natl. Acad. Sci. USA* **92** 5465-5469
- Phillips K L, Gartrell D M, Roses A D and Lee J E 1993 A triplet repeat polymorphism in a gene expressed in human hypothalamus; *Hum. Mol. Genet.* **2** 1332
- Polymeropoulos M H, Rath D S, Xiao H and Merrill CR 1991 Trinucleotide repeat polymorphism at the human transcription factor IID gene; *Nucleic Acids Res.* **19** 4307
- Richards R I, Holman K, Yu S and Sutherland G R 1993 Fragile X syndrome unstable element. p(CCG)_n, and other simple tandem repeat sequences are binding sites for specific nuclear proteins; *Hum. Mol. Genet.* **2** 1429-1435
- Richards R I and Sutherland G R 1992 Dynamic mutations: A new class of mutations causing human disease; *Cell* **70** 709-712
- Riggins G J, Lokey L K, Chastain J L, Leiner H A, Sherman S L, Wilkinson K D and Warren S T 1992 Human genes containing polymorphic trinucleotide repeats; *Nature Genet.* **2** 186-191
- Stallings R L 1994 Distribution of trinucleotide microsatellites in different categories of mammalian genomic sequence: Implications for human genetic diseases; *Genomics* **21** 116-121
- Streisinger G, Okada Y, Emnch J, Newton J, Tsugita A, Terzaghi F and Inouye M 1966 Frameshift mutations in the genetic code; *Cold Spring Harb. Symp. Quant. Biol.* **31** 77-84
- Strand M, Prolla T A, Liskay R M and Petes T D 1993 Destabilisation of tracts of simple repetitive DNA in yeast by mutations affecting DNA mismatch repair; *Nature (London)* **365** 274-276
- Stott K, Blackburn J M, Butler P J G and Perutz M 1995 Incorporation of glutamine repeats makes protein oligomerize: Implications for neurodegenerative diseases; *Proc. Natl. Acad. Sci. USA* **92** 6509-6513
- Sutherland G R and Richards R I 1995 Simple tandem DNA repeats and human genetic disease; *Proc. Natl. Acad. Sci. USA* **92** 3636-3641
- Trifonov E N 1989 The multiple codes of nucleotide sequences; *Bull. Math. Biol.* **51** 417-432
- Tripathi J and Brahmachari S K 1991 Distribution of simple repetitive (TG/CA)_n and (CT/AG)_n sequences in human and rodent genomes; *J. Biomol. Struct. Dyn.* **9** 387-397
- Wang Y H, Amirhaeri S, Kang S, Wells R D and Griffith J D 1994 Preferential nucleosome assembly at DNA triplet repeats from the myotonic dystrophy gene; *Science* **265** 669-671