

Trinucleotide repeats and long homopeptides in genes and proteins associated with nervous system disease and development

SAMUEL KARLIN AND CHRIS BURGE

Department of Mathematics, Stanford University, Stanford, CA 94305-2125

Contributed by Samuel Karlin, October 27, 1995

ABSTRACT Several human neurological disorders are associated with proteins containing abnormally long runs of glutamine residues. Strikingly, most of these proteins contain two or more additional long runs of amino acids other than glutamine. We screened the current human, mouse, *Drosophila*, yeast, and *Escherichia coli* protein sequence data bases and identified all proteins containing multiple long homopeptides. This search found multiple long homopeptides in about 12% of *Drosophila* proteins but in only about 1.7% of human, mouse, and yeast proteins and none among *E. coli* proteins. Most of these sequences show other unusual sequence features, including multiple charge clusters and excessive counts of homopeptides of length \geq two amino acid residues. Intriguingly, a large majority of the identified *Drosophila* proteins are essential developmental proteins and, in particular, most play a role in central nervous system development. Almost half of the human and mouse proteins identified are homeotic homologs. The role of long homopeptides in fine-tuning protein conformation for multiple functional activities is discussed. The relative contributions of strand slippage and of dynamic mutation are also addressed. Several new experiments are proposed.

Several human dominantly inherited neurodegenerative diseases are associated with abnormally expanded tracts of glutamine codons (for recent reviews, see refs. 1–3). Principal examples include Huntington disease (HD) protein (see Table 1), containing near the N terminus Q_{23} (Q = glutamine residue in the one-letter code); ataxin I (spinocerebellar ataxia type 1, SCA1), containing the glutamine cluster $Q_{12}HQHQ_{15}$; atrophin 1 (denatorubral–pallidolusian atrophy 1, DRPLA), containing the homopeptide Q_{20} ; and androgen receptor (AR; spinobulbar muscular atrophy, Kennedy disease), containing the glutamine homopeptides Q_{21} , Q_6 , Q_5 . In these proteins the length of the polyglutamine run is polymorphic but of limited size in healthy individuals (1). Disease severity appears to be correlated with the extent of tandem reiteration of the CAG (glutamine) codon above a certain threshold number (1, 2).

Proposed consequences of hyperexpansion of DNA triplet repeats (especially of CTG and CAG) include loss of or altered rate of transcription or translation, mRNA instability, aberrant DNA hairpin structures, etc. (1–3). Hypermethylation has been shown as a consequence of extended reiterations of CGG (1, 2). On the protein structural level, difficulties may arise from protein aggregation due to attachment of glutamine-rich proteins to unrelated molecules (6), to inappropriate multimerization, or to formation of “polar zippers,” in which a long stretch of glutamine residues link β -strands by hydrogen bonds (7–9).

It is striking that three of the trinucleotide disease proteins contain multiple long homopeptides, not only of glutamine (Table 1). Specifically, listing all homopeptides of lengths \geq five residues, HD also contains P_{11} , P_{10} , E_5 , and E_6 ; DRPLA

also contains S_7 , S_{10} , P_6 , and H_5 ; and AR also contains P_8 , A_5 , and G_{24} . What kinds of proteins contain multiple long homopeptides? What specific functions or structures are associated with particular long homopeptides? Under what conditions will hyperexpansion of a DNA triplet occur, and when will such expansions cause disease? Why are the triplet-repeat disease genes conspicuous in neuronal tissues?

We identified all human, mouse *Drosophila melanogaster*, yeast (*Saccharomyces cerevisiae*), and *Escherichia coli* protein sequences of length \geq 200 residues that contain multiple long homopeptides. (For precise criteria and rationale, see the legend to Table 1.) A search of the SWISS-PROT data base release 31 (February 1995) for sequences containing multiple long homopeptides identified 36 of 2234 (1.6%) human sequences (Table 1), 68 of 551 (12%) *D. melanogaster* sequences (Table 2), 22 of 1242 (1.8%) mouse sequences (not shown), 42 of 2454 (1.7%) yeast sequences (not shown), and 0 of 1790 *E. coli* sequences. These examples feature homopeptides predominantly of the uncharged polar amino acids Q, N, S, T, P, and H; the small residues G and A; or acidic types D and E. Apart from three cases of leucine runs among the human examples, no other aliphatic or aromatic runs appeared.

The examples of Tables 1 and 2 are intriguing. (i) With few exceptions, the *Drosophila* proteins contribute prominently in developmental regulatory capacities, and the majority of these proteins are active in the central or peripheral nervous system (CNS or PNS). (ii) Most of these proteins contain many unusual sequence features, including significantly high multiplet counts, multiple charge clusters, and multiple significantly long uncharged segments (see below for definitions and further discussion). (iii) Several contrasts between the *Drosophila* examples (Table 2) and the human examples (Table 1) stand out. (a) The proportion of *Drosophila* proteins with multiple long homopeptides (12%) is dramatically higher than in human proteins (1.6%). (b) In the *Drosophila* sequences, the predominant homopeptides are glutamine runs, found in 51 of 68 cases, whereas polyglutamine (Q_n , $n \geq 5$) occurs in only six of the 34 human proteins. (c) The *Drosophila* proteins tend to be longer than the human proteins (median amino acid length of sequences: human, 527; *Drosophila*, 748; and mouse, 460).

Before discussing in detail the sequences of Tables 1 and 2, we present some background discussion on charge clusters, uncharged runs, and multiplet counts in protein sequences.

BACKGROUND AND METHODS

Charge Clusters and Uncharged Runs. A charge cluster is a protein segment (typically 20–75 residues) with significantly high specific charge content relative to the charge composition of the whole protein; see refs. 5, 10, 11 for precise characterizations. The percentage of proteins with at least one significant charge cluster is about 20–25% in most eukaryotic species (e.g., human, mouse, and yeast), about 35% in *Dros-*

Table 1. Sequence features of human triplet repeat disease proteins and other human proteins containing multiple long homo-amino acid runs

Protein name (accession no.)	Length, aa	Homo-amino-acid runs of length ≥ 5 aa	Charge clusters +, -, \pm	Uncharged runs, no.	Multiplets (aa, codon)
Triplet repeat disease proteins					
Androgen Receptor (Kennedy's Disease) (P10275)	919	Q ₂₁ , Q ₆ , Q ₅ , P ₈ , A ₅ , G ₂₄	1, 0, 0	0	0, +
Atrophin 1 (DRPLA) (N/A)	1191	S ₇ , S ₁₀ , P ₆ , H ₅ , Q ₂₀	0, 0, 3	5	0, n/a
Ataxin 1 (SCA1) (S46268)	816	Q ₁₂ , Q ₁₅	0, 0, 2	2	0, 0
Huntingtin (Huntington Disease) (P42858)	3144	Q ₂₃ , P ₁₁ , P ₁₀ , E ₅ , E ₆	0, 1, 0	1	++, ++
Other human proteins with multiple long homopeptides					
RDC-1 (Q01851)	331	H ₉ , G ₅ , G ₆ , G ₆ , A ₅ , A ₁₀	0, 0, 1	1	0, 0
Nervous system-specific Octamer Binding Factor N-OCT 3 (P20265)	443	G ₂₁ , Q ₂₁ , A ₅ , P ₇	0, 0, 2	1	+, ++
Calcium-dependent Protease, Small Subunit (P04632)	268	G ₁₁ , G ₂₀	0, 0, 0	2	0, 0
SOX-3 Protein (P41225)	443	G ₅ , A ₁₅ , P ₅ , A ₇ , A ₆ , A ₁₀	1, 0, 1	1	0, 0
Sarcoplasmic Reticulum Histidine-rich Calcium-binding Protein (P23327)	699	E ₁₂ , D ₁₆ , E ₇ , H ₅ , E ₆ , E ₇ , E ₈	0, 2, 0	1	0, +
Thyroliberin (P20396)	242	E ₆ , E ₇ , E ₈	0, 0, 0	1	0, n/a
Hox-B4 (P17483)	251	P ₁₅ , P ₅	0, 0, 0	0	0, n/a
Octamer-binding Transcription Factor 6 (Q03052)	398	G ₆ , G ₆ , G ₆ , P ₆ , H ₆	0, 0, 2	1	0, 0
OTX1 (P32242)	354	S ₅ , A ₆ , H ₅ , H ₁₀	1, 0, 1	3	+, n/a
SOX-4 (Q06945)	474	G ₇ , G ₆ , G ₆ , A ₇ , S ₇	0, 1, 1	0	+, ++
Engrailed-1 (Q05925)	391	P ₆ , A ₁₀ , A ₉	0, 0, 1	1	0, ++
Factor VIII Intron 22 Protein (P23610)	364	A ₅ , A ₅ , P ₆ , L ₆	0, 0, 0	0	0, 0
Major Centromere Autoantigen B (P07199)	599	E ₁₄ , E ₆ , E ₆ , E ₅ , D ₅	0, 1, 0	0	++, ++
Translational Initiation Factor 2 β (P20042)	333	K ₈ , K ₆ , K ₆	0, 0, 0	0	0, 0
Heterogeneous Nuclear Ribo-NP L (P14866)	558	G ₇ , G ₉ , G ₆ , P ₈	0, 0, 0	0	0, 0
Acrosin (P10323)	421	P ₅ , P ₈ , P ₉	0, 0, 0	0	0, 0
Transcriptional Repressor Prot. YY1 (P25490)	414	E ₅ , G ₅ , H ₁₁	0, 1, 0	1	++, ++
Hox-B3 Protein (P14651)	431	G ₁₁ , G ₁₀	1, 0, 1	0	0, 0
Chromogranin A (P10645)	457	E ₈ , E ₉ , E ₅	0, 0, 0	1	0, 0
Protein Kinase C Substrate, 80 KD (P14314)	527	L ₅ , E ₁₂ , E ₇	0, 1, 0	0	0, 0
Ring3 (P25440)	754	E ₆ , E ₅ , E ₅ , S ₅ , S ₁₁	1, 1, 1	2	++, +
Early Growth Response 2 (P11161)	406	P ₇ , A ₁₀	0, 0, 1	1	0, 0
Early Growth Response 1 (P18146)	543	S ₆ , G ₉ , S ₆	1, 0, 0	1	0, 0
Nucleolin (P19338)	706	A ₅ , D ₅ , D ₅ , E ₁₀	0, 2, 1	0	++, ++
Amyloid-like Protein 2 (Q06481)	763	A ₅ , L ₆ , E ₅ , E ₇	0, 1, 0	2	0, 0
MHC Class II Reg. Factor RFX-1 (P22670)	979	P ₇ , G ₁₄ , E ₈	1, 1, 0	2	0, 0
Insulin Receptor Substrate-1 (P35568)	1242	G ₅ , S ₇ , Q ₆ , G ₅ , S ₅ , P ₇	0, 0, 0	0	++, +
SP4 Transcription Factor (Q02446)	784	E ₅ , A ₆ , S ₉	0, 1, 1	2	0, n/a
Heterogeneous Nuclear Ribo-NP U (Q00839)	806	E ₁₁ , G ₆	0, 1, 0	0	++, ++
Macrophage CSF I Receptor (P07333)	972	L ₅ , L ₁₀	0, 0, 0	0	0, 0
HRX (Q03164)	3969	G ₇ , A ₇ , S ₉ , G ₅ , E ₅	0, 0, 1	6	++, ++
Adenomatous Polyposis Coli (P25054)	2843	S ₅ , P ₅ , A ₅ , S ₆	0, 2, 0	0	0, 0

SWISSPROT release 31 (Feb 1995) was screened for human proteins satisfying either or both of the following: (i) three or more homopeptides of length ≥ 5 amino acids (aa), whose combined lengths total at least 20 aa; or (ii) at least one homopeptide of length ≥ 10 aa and at least one other of length ≥ 5 aa. Proteins with extremely biased aa composition (1 or more aa with frequency $> 20\%$) were excluded, since long homopeptides are expected to occur quite frequently by chance in sequences of very biased composition. For each group of two or more substantially identical proteins, only the longest was chosen for analysis to avoid redundancy. This resulted in 36 out of 2234 (1.6%) human proteins being chosen, including the three triplet-repeat disease gene products androgen receptor, ataxin 1, and huntingtin. The atrophin 1 (DRPLA) sequence, not in SWISSPROT but included in the counts above for simplicity, was provided to us by R. Margolis (Johns Hopkins Medical School). The human TATA box-binding protein TFIID, which contains Q₃₆ but no other long homopeptide, was excluded by these criteria. The proteins in each of the groups (disease-related and other proteins) were ordered on the basis of their multiplicity index, $M = (\text{no. of aa in runs of length } \geq 5 \text{ aa}) / (\text{protein length})$, from highest to lowest. Column 2 lists all homopeptide runs of length ≥ 5 aa in each protein in the order in which they occur along the protein sequence. Column 3 lists the number of significant positive charge clusters (clusters of R and K, one letter code), negative charge clusters (D and E), and mixed charge clusters (R, K, D, and E). Charge clusters are defined statistically as regions unusually rich in a particular charged residue type (4). Column 4 lists the number of significantly long uncharged runs—unusually long stretches devoid of charge or with only one or two charged residues (4, 5). Column 5 indicates the significance of the number of amino acid “multiplets” and codon “multiplets”—patterns of the type X_n , where X is any of the 20 amino acid types or one of the 61 codon types, respectively, and $n \geq 2$ (runs longer than 2 are counted only once). The number of multiplets in a protein (or its corresponding gene sequence) is evaluated in terms of the number of standard deviation (SD) units above or below the mean, given the length and amino acid/codon composition of the protein/gene (4): counts > 2.3 SD above the mean are indicated by “+”; those > 3.0 SD above the mean, by “++”; those within 2.3 SD of the mean, by “0.” (No negative extremes were observed.) n/a, gene not available; Dis., disease.

ophila, and only 7% in *E. coli* (and most prokaryotes) (10). In eukaryotes, charge clusters are associated with transcriptional activation, membrane receptor activity, and developmental regulation, the last category apparently overrepresented among available *Drosophila* sequences (see refs. 10 and 11 for specific examples and classifications). In contrast, hardly any

distinctive charge configurations occur among the bulk of housekeeping proteins, cytoplasmic enzymes, or among prokaryotic proteins (10). Proteins with multiple charge clusters are uncommon, about 3.5% in human, mouse, and yeast, but are relatively more frequent (13.5%) in *Drosophila* (11). Primary families of proteins with multiple charge clusters are

Table 2. Sequence features of *Drosophila* proteins containing multiple long homo-amino acid runs

Protein	Length, aa	Homo-amino acid runs of length ≥ 5 aa	Charge clusters +, -, \pm	Uncharged runs, no.	Multiplets (aa, codon)	Devel./Neuro./TF	Description/function
E74-B (P11536)	883	S ₁₄ , S ₈ , A ₅ , Q ₇ , Q ₈ , A ₁₃ , Q ₈ , Q ₅ , Q ₁₁ , A ₉ , S ₅ , A ₅	2, 2, 3	3	++, ++	D, T	Control of molting, ecdysome inducible
Abdominal-B (P09087)	491	Q ₁₀ , Q ₈ , Q ₆ , Q ₅ , Q ₁₀ , Q ₇ , A ₆ , N ₅	1, -, 1	1	+, ++	D, N, T	Segment/PNS devel.
Odd-skipped (P23803)	392	Q ₁₂ , Q ₆ , Q ₈ , H ₆ , S ₇	0, 1, 1	1	0, ++	D, N, T	Pair-rule segmentation
Big brain (P23645)	700	Q ₁₃ , Q ₁₀ , P ₅ , Q ₂₅ , Q ₅	0, 1, 1	1	0, n/a	D, N	Neural/epidermal devel.
Hairy (P14003)	337	Q ₅ , Q ₆ , Q ₆ , A ₁₀	0, 0, 1	2	0, 0	D, N, T	Segment/sense organ/PNS development
Dorsal (P15330)	678	Q ₅ , Q ₈ , Q ₆ , Q ₁₄ , A ₅ , Q ₇	0, 0, 1	2	0, ++	D	Controls dorsal polarity
AEF1 (P39413)	308	N ₅ , Q ₁₃ , Q ₆	0, 0, 0	2	+, n/a	T	Regulates <i>Adh</i>
E75-C (P13055)	1443	Q ₅ , Q ₇ , P ₁₅ , Q ₆ , Q ₆ , S ₇ , S ₆ , Q ₈ , P ₈ , S ₅ , S ₅ , Q ₉ , Q ₅ , S ₁₇	1, 0, 0	2	++, ++	D, T	Control of molting, ecdysome inducible
Zeste (P09956)	574	Q ₆ , Q ₇ , A ₆ , Q ₅ , Q ₇ , A ₆ , G ₅	1, 1, 2	1	0, ++	D, T	Activates <i>Ubx</i>
Ubx (P02834)	389	A ₅ , G ₁₃ , A ₁₀	1, 0, 1	2	0, n/a	D, T	Anter./post. devel.
Giant (P39572)	448	Q ₁₁ , Q ₁₂ , A ₉	1, 0, 0	1	++, ++	D, T	Represses Kruppel, Knirps
Caudal (P09085)	472	H ₅ , N ₁₀ , N ₇ , R ₁₁	1, 0, 1	2	0, n/a	D, T	Segment devel. (early)
Abdominal-A (P29555)	330	Q ₁₇ , Q ₆	0, 0, 1	1	0, +	D, N, T	Segment/PNS devel.
Engrailed (P02836)	552	Q ₁₁ , A ₁₄ , A ₆ , S ₇	0, 1, 1	3	++, ++	D, N, T	Intrasegment devel.
P62/RM62 (P19109)	575	G ₆ , G ₇ , G ₅ , G ₆ , G ₇ , G ₈	0, 0, 0	1	0, 0	—	RNA helicase
C/EBP (Q02637)	444	Q ₂₀ , G ₅ , Q ₅	1, 1, 0	1	++, n/a	D, T	Border cell migration
Prospero (P29617)	1407	A ₆ , A ₅ , Q ₆ , N ₇ , N ₅ , D ₅ , Q ₁₈ , Q ₆ , Q ₅ , Q ₂₄ , P ₇	0, 2, 1	4	++, ++	D, N, T	Embryonic CNS devel.
Female sterile homeotic (P13709)	2038	A ₁₄ , S ₇ , S ₇ , Q ₁₀ , Q ₁₇ , Q ₆ , Q ₆ , Q ₅ , Q ₆ , Q ₁₁ , A ₆ , S ₇ , S ₅ , A ₆ , G ₅ , G ₈ , G ₆	0, 1, 5	6	++, ++	D	Oogenesis, regulates homeotic genes
RBP-J kappa (P28159)	594	Q ₆ , Q ₆ , Q ₁₄ , N ₆ , Q ₆	1, 0, 1	2	0, ++	D, N, T	Suppressor of hairless
Maternal Pumilio (P25822)	1533	A ₆ , G ₇ , A ₅ , Q ₈ , Q ₁₃ , A ₁₁ , Q ₁₅ , A ₁₁ , Q ₆ , A ₁₀ , A ₆	1, 0, 1	4	++, ++	D	Anter./post. devel.
Ftz-F1 (P33244)	1043	Q ₉ , A ₆ , T ₅ , N ₇ , N ₇ , T ₅ , Q ₇ , G ₉ , G ₅ , G ₆	1, 1, 0	3	++, ++	D, T	Activation of <i>Ftz</i>
DSX, male (P23023)	549	A ₈ , A ₆ , H ₅ , G ₉ , P ₆	1, 1, 2	2	++, n/a	D, T	Controls sexual diff.
ELAV (P16914)	483	A ₆ , A ₆ , Q ₅ , Q ₆ , A ₆	0, 0, 1	1	0, 0	D, N	Nerve development
Deformed (P07548)	590	G ₆ , Q ₅ , N ₅ , N ₁₁ , N ₇	1, 1, 1	3	++, ++	D, N, T	Regulates segmentation
DGK (Q01583)	517	S ₅ , A ₅ , Q ₁₃ , Q ₆	0, 0, 0	0	0, 0	—	Signal transduction
Knirps (P10734)	429	A ₅ , Q ₇ , A ₆ , S ₆	1, 0, 0	0	++, ++	D, T	Abdominal segmentation
Orthodenticle (P22810)	671	Q ₅ , A ₅ , A ₅ , A ₅ , A ₇ , G ₁₀	1, -, 1	1	++, ++	D, N, T	Segmentation (head)
Topoisomerase I (P30189)	972	H ₆ , S ₇ , S ₅ , S ₆ , S ₅ , S ₅ , S ₈ , S ₆ , D ₅	0, 1, 1	2	0, ++	—	Alters DNA supercoiling
Eyes absent (Q05201)	766	Q ₇ , Q ₉ , G ₅ , A ₆ , A ₇ , A ₇	2, 1, 1	1	++, ++	D, T	Eye development
Antennapedia (P02833)	378	Q ₉ , Q ₅ , Q ₆	1, 0, 1	2	0, 0	D, T	Anter./post. devel.
Cut (P10180)	2175	A ₅ , A ₉ , Q ₇ , A ₉ , N ₇ , N ₆ , N ₅ , A ₁₁ , Q ₅ , A ₅ , A ₇ , A ₅ , A ₇ , A ₅ , A ₈ , P ₆ , S ₆	0, 2, 2	5	++, ++	D, N, T	Embryonic CNS devel.
Hairless (Q02308)	1077	N ₅ , T ₅ , S ₅ , S ₅ , S ₇ , A ₉ , A ₁₀ , A ₈	1, 1, 1	2	++, ++	D, N, T	Sense organ/PNS devel.
Polycomb (P26017)	390	H ₁₀ , H ₈	0, 0, 0	0	0, ++	D, N	Regulates homeotic genes
Polyhomeotic-proximal chromatin (P39769)	1589	Q ₇ , Q ₅ , Q ₆ , Q ₅ , Q ₆ , Q ₇ , Q ₆ , Q ₅ , Q ₇ , Q ₇ , Q ₆ , T ₆	1, 0, 2	3	++, n/a	D, T	Regulates segment development
Disconnected (P23792)	568	A ₆ , P ₅ , A ₈ , H ₇	1, 2, 1	3	++, ++	D, N, T	Embryonic nerve devel.
5HT-receptor (P28285)	834	N ₅ , T ₅ , A ₅ , A ₅ , G ₅ , A ₅ , Q ₈	1, 2, 0	1	0, +	N	Serotonin receptor
Shaker- ϵ (P08513)	656	Q ₆ , Q ₇ , Q ₁₁ , Q ₅	1, 0, 2	1	++, ++	D, N	Voltage-dep. K ⁺ channel
Labial (P10105)	635	Q ₅ , Q ₅ , Q ₆ , Q ₆ , G ₆	-, -, 1	0	++, ++	D, T	Anter./post. devel.
GRH/ELF1 (P13002)	1063	Q ₉ , A ₈ , Q ₁₇ , Q ₅ , Q ₇	0, 0, 1	3	+, +	D, T	Nerve devel./cuticle
Suppressor 2 of zeste (P25172)	1365	Q ₁₁ , Q ₁₄ , Q ₅ , N ₂₀ , S ₉	0, 1, 1	0	0, 0	D, N, T	Segmentation, regulates <i>Ubx</i>
Yan (Q01842)	732	Q ₁₀ , Q ₈ , Q ₅ , A ₇	0, 1, 1	2	0, ++	D, N, T	Neural differentiation
Daughterless (P11420)	710	G ₅ , Q ₅ , G ₇ , A ₅ , Q ₇	0, 0, 1	2	0, ++	D, N, T	Neurogenesis, sex determ.
Stubble (Q05319)	786	S ₁₀ , Q ₆ , T ₈ , S ₈	0, 1, 0	0	+, +	D	Epithelial devel.
Empty spiracles (P18488)	497	Q ₆ , Q ₆ , A ₈	1, 1, 1	1	+, ++	D, T	Antenna/mandible devel.
KNRL (P13054)	647	A ₆ , G ₉ , N ₁₁	1, 2, 0	2	++, ++	D, T	Embryonic segmentation
PSC (P35820)	1603	T ₅ , T ₆ , T ₈ , T ₆ , S ₁₇ , S ₁₀ , P ₅	0, 1, 1	1	++, ++	D, N, T	Regulates homeotic genes
Runt (P22814)	509	A ₁₂ , A ₆	0, 1, 0	2	0, +	D, N, T	Pair-rule segmentation
Fork head (P14734)	510	G ₁₃ , H ₅	0, 0, 1	1	+, ++	D, T	Embryonic CNS devel.
CF2, III (Q01522)	514	P ₆ , Q ₁₁	0, 0, 0	0	0, 0	D, T	Chorion development
SHAB11 (P17970)	924	Q ₇ , Q ₉ , Q ₅ , Q ₉	0, 0, 2	3	+, ++	N	Voltage-dep. K ⁺ channel
TKR (P14083)	753	E ₅ , A ₆ , A ₆ , A ₈	0, 1, 0	1	++, n/a	D	Protein-tyrosine kinase
SOL (P27398)	1597	Q ₆ , H ₅ , H ₆ , S ₅ , Q ₅ , Q ₉ , H ₈ , G ₅	0, 0, 0	0	0, ++	N	Neuronal devel.

Table 2. (Continued)

Protein Name (accession no.)	Length, aa	Homo-amino acid runs of length ≥ 5 aa	Charge			Devel./ Neuro./ TF	Description/function
			clusters +, -, \pm	Uncharged runs, no.	Multiplets (aa, codon)		
Rutabaga (P32870)	2248	G ₁₀ , Q ₅ , Q ₈ , Q ₅ , Q ₅ , Q ₅ , H ₆ , Q ₅ , D ₅ , Q ₁₃	1,2,1	4	++, ++	N	Mutant has learning defect
Single-minded (P05709)	673	S ₅ , Q ₅ , S ₅ , Q ₅	0,0,1	1	+, +	D, N, T	CNS development
Staufen (P25159)	1026	Q ₅ , P ₅ , Q ₇ , P ₅	0,0,1	2	++, ++	D	Controls oocyte polarity
Inflated (P12080)	1394	S ₇ , S ₅ , S ₆ , S ₁₀	0,0,0	2	++, 0	D	Intercellular comm.
PTP99A (P35832)	1301	S ₅ , Q ₁₆ , S ₅	0,0,0	2	+, ++	D, N	Intercellular comm.
Zinc-finger 1 (P28166)	1060	Q ₆ , Q ₅ , Q ₅ , Q ₅	0,2,1	2	+, +	D, N, T	Embryonic CNS devel.
DWNT-3 (P28466)	1010	G ₆ , S ₁₄	0,0,0	0	0, +	D, N	CNS/wing devel.
Spalt-major (P39770)	1356	S ₅ , A ₆ , A ₁₄	0,2,2	4	++, n/a	D, N, T	Segment development
Periodic clock (P07663)	1155	G ₅ , A ₅ , A ₆ , A ₅	0,1,1	3	+, +	N	Biological rhythms
SU(S) (P22293)	1322	G ₅ , G ₅ , Q ₆ , G ₈	1,2,1	2	++, ++	D	RNA-binding suppressor
Notch (P07207)	2703	A ₈ , G ₉ , Q ₁₃ , Q ₁₇	0,0,0	2	0, 0	D, N	Neurogenic, oogenic
Brahma (P25439)	1638	Q ₅ , Q ₆ , E ₅ , E ₅ , D ₆	0,3,2	4	++, ++	D, T	Regulates homeotic genes
SOS (P26675)	1595	G ₅ , G ₅ , A ₆ , Q ₆	0,0,0	1	0, 0	D, N	Neuronal development
Trithorax (P20659)	3759	S ₅ , D ₆ , Q ₇ , Q ₅ , Q ₇ , S ₇	0,2,1	6	++, ++	D, T	Segmentation
Zinc-finger 2 (P28167)	3005	S ₅ , A ₇ , Q ₅ , S ₇ , Q ₅	0,0,0	4	++, 0	D, N, T	Embryonic CNS devel.
FAT (P33450)	5147	S ₅ , D ₅ , S ₇ , Q ₅ , Q ₆ , P ₅ , Q ₆	0,0,1	0	0, 0	D, N	Cell adhesion protein

SWISSPROT release 31 (February 1995) was screened for *D. melanogaster* proteins satisfying the criteria described in the legend to Table 1. This resulted in 68 out of 551 *Drosophila* proteins (12%) being chosen. The layout of Table 2 is the same as that of Table 1, except that two new columns have been added. Column 7 (Devel./Neuro./TF) describes whether or not the protein is a developmental protein (D), a nervous system-related protein (N), and/or a transcription factor (T). Column 8 provides a very brief indication of the major developmental role(s) or other known functions of each protein. A "-" in column 4 indicates that the frequency of the corresponding charged residue type was too low in the given protein for evaluation of charge clusters. The major *Drosophila* developmental protein mastermind, which exhibits (in order of occurrence) the plethora of runs Q₁₀, Q₅, Q₁₆, Q₅, N₅, G₁₀, N₅, G₇, N₅, Q₅, Q₇, Q₁₄, Q₆, Q₇, Q₆, Q₁₆, Q₅, Q₅, Q₇, Q₈, Q₆, Q₅, Q₆, A₁₀, Q₅, Q₁₂, Q₅, T₅, G₈, and G₆, is excluded because of its biased composition (Q = 21.7%). diff., Differentiation; anter./post., anterior/posterior; dep., dependent; determ., determination; comm., communication; 5HT, serotonin; CNS and PNS, central and peripheral nervous systems.

(i) developmental regulatory proteins, (ii) voltage-gated ion-channel proteins, and (iii) regulatory proteins of large eukaryotic DNA viruses (10, 11).

Sequence Features of Transcription Factors. Transcription factors tend to have a modular organization, generally involving a DNA-binding domain, activating domain(s), and other structural/functional domains separated by linker segments (12, 13). Generally, transactivating domains in transcription factors have been associated with acidic clusters (e.g., GAL-4, GCN-4), with proline or glutamine clusters (e.g., SP1, CREB), or with serine/threonine-rich regions (e.g., POU proteins) (for reviews, see refs. 12 and 13). For example, the charge cluster overlapping the POU domain of human Oct-1 is flanked by a Q-rich region on the N-terminal side and by an abundance of S and T residues on the C-terminal side. The human transcription factor Sp1 has a charge cluster centering on three tandem zinc-finger motifs that are responsible for DNA binding and a long uncharged region near the N terminus rich in S, T, and Q residues. The Q-rich portions are considered important for transcriptional activation (14, 15). Recently, based on a series of GAL-4 chimeric constructs manipulating different sizes of Q_n-runs, the highest levels of transcription were achieved with polyglutamine tracts in the length range of 6 to 34 residues (16). Longer tracts showed a reduced ability to mediate transcriptional activation. Similar results were achieved for P_n runs, but with a steeper dropoff in activation activity for longer runs (16). Multiple significant charge clusters, in conjunction with one or more significantly long uncharged stretches of mostly polar residues, occur in many *Drosophila* developmental regulatory proteins (e.g., engrailed, paired, cut, deformed, and doublesex; see Table 2). Among the trinucleotide disease genes, only the androgen receptor is an established transcription factor.

Multiplet Counts. A measure of the homopeptide density of a protein sequence is the multiplet count: the number of distinct homooligopeptide runs of 2 or more residues (for details, see the legend to Table 1 and ref. 4). A complete absence of proteins with significantly high multiplet counts in

E. coli and a pronounced overrepresentation of proteins with significant multiplet counts in *Drosophila* (about 10%) have been observed (4). The percentage of human proteins with significant multiplet counts is about 1.6% with similar percentages observed in mouse and yeast. We also assessed for the genes of Tables 1 and 2 the multiplet count on the codon level, indicating the number of distinct homocodon runs of length two or more. Significance is evaluated similarly to the amino acid case.

RESULTS

Distinctive Sequence Features of Proteins with Multiple Long Homopeptides. The presence of multiple long homopeptides correlates strongly with the presence of unusual sequence features including multiple charge clusters, significantly long uncharged segments, and high multiplet counts (Tables 1 and 2). It is useful to highlight and contrast the distinctive sequence features of the proteins with multiple long homopeptides versus general proteins.

D. melanogaster (68 with multiple long homopeptides). (i) Fifty-five (80%) are involved in developmental control, and 31 (45%) are involved in CNS development; (ii) 47 (69%) have multiple distinct charge clusters (overall *Drosophila*, 13.5%); (iii) only 12 (18%) are absent of any charge cluster (overall, 65%); (iv) 38 (56%) have significant amino acid multiplet counts (overall, 12%); 40 of 63 (five gene sequences were unavailable) show significant codon multiplet counts, including most of those with high counts of amino acid multiplets.

Human (36 with multiple long homopeptides). (i) Eight (22%) have at least two charge clusters (overall, 3–4%); (ii) 25 (69%) have at least one charge cluster (overall, 20–25%); (iii) 8 (22%) have significant amino acid multiplet counts (overall, <1.5%); and HOX-B3, HOX-B4, RDC-1, engrailed-1, N-OCT3, OTX-1, and OCT6 are likely developmental proteins (N-OCT3 participates in neural development).

Mouse (22 with multiple long homopeptides). (i) Eight (36%) have multiple charge clusters; (ii) 18 (82%) have at least one

charge cluster; (iii) 5 (23%) have significant multiplet count (including the mouse homolog of AR); apparently active in development are BRN1 (brain specific homeobox), HME1 (engrailed-1), HOX-AP, HOX-B4, HOX-D8, HOX-D9, HOX-D11, OTX1, FORMIN-4, N-OCT3, and OCT6.

Long acidic runs. There are several examples of Table 1 that possess hyperacidic runs characterized in ref. 11: CENP-B, nucleolin, sarcoplasmic reticulum, histidine-rich calcium-binding protein (HRC), chromogranin A, thyroliberin, and amyloid-like protein 2 (APLP2). Activity of several of these proteins has been detected in association with chronic systemic autoimmune diseases (17). Actually, more than half of all proteins eliciting autoimmune diseases show unusual charge configurations and especially long acidic runs (11, 17), suggesting that hyper-expansion of acidic homopeptides might lead to illness. The listing of Table 1 also includes the oncogenes adenomatous polyposis coli (APC) protein and HRX (chromosomal translocation associated with leukemia).

Codon content of long homopeptides. The glutamine runs of the triplet-repeat disease proteins are almost entirely encoded from the codon CAG. However, most of the other long homopeptides show considerable variation in codon usage. For example, P₈ in the androgen receptor is coded by (CCG)₂(CCC)-(CCT)(CCG)₂(CCT)(CCC). Interestingly, G₂₄ is encoded from (GGT)₃(GGG)(GGT)(GGC)₁₈. The P₁₁ peptide of the HD protein immediately downstream of Q₂₃ is encoded from (CCG)(CCA)(CCG)₇(CCT)₂. The nervous system-specific octamer binding factor N-OCT3 has Q₂₁ encoded by a genuine mixture of CAG and CAA codons, and the same is true for most examples of Tables 1 and 2. Apart from a few serine runs encoded exclusively from AGC—e.g., early growth response protein-1 (EGR1) and insulin receptor substrate-1 (IRS1)—most other serine runs of Table 2 use both serine codon types AGY and TCN. In general, the codon content of long homopeptides does not present a consistent pattern. Homopeptides encoded by iterations of a single codon are likely the result of strand slippage; substantial use of multiple codons suggests possible importance of the homopeptide for protein structure or function.

In addition to the multiple long homopeptides, some of the human triplet repeat disease genes contain unusual alternating patterns.

Histidine patterns. The alternating histidine pattern (HX)₈ = H₂HQHSIHSHLHLHQ in the DRPLA protein (Table 1) is notable in view of the plethora of histidine periodic patterns in *Drosophila* developmental proteins. These include the *Drosophila* female sterile homeotic protein (FSH; see Table 2), featuring (HG)₄¹ = HH(HG)₁₃ (the superscript indicates the number of errors—i.e., departures from the pattern); homeotic bicoid, (HX)₁₂¹ = (HT)₂(HP)₂HS(HP)₂HS(HP)₂HHQHP; segmentation protein paired, (HX)₈²; homeotic spalt, [(E/G)HH]₇; homeobox protein NK-1, (HX)₁₃²; male specific DsX showing three separated patterns, (HHX)₈², (HX)₆¹, and (HX)₈²; ecdysone-inducible proteins E74A, (HX)₁₁¹, and E75B, (HQ)₅; segmental pattern protein deformed, (HX)₁₀² and (HX)₆; protein grainy-head (ELF1), (HX)₈¹; homeobox protein HM(10), (HX)₂₃². *Drosophila* extracellular structural proteins chorion 36 and chorion 38 show (GH)₅¹ and (GH)₁₁², respectively. There are also many developmental proteins with long histidine runs—e.g., polycomb H₁₀ and H₈ and hunchback H₆. Histidine is a versatile amino acid that can adopt flexible roles in conformation, in catalytic actions, in metal coordination, and various enzymatic activities. Histidine patterns and runs provide opportunities for differential charge gradients, hydrogen bonding, and metal coordination. The only mouse proteins with long histidine runs found in the current SWISSPROT data base are the homeotic protein HXA1, featuring H₁₁; the homeobox proteins MOX-2A, containing H₁₂; and the regulatory nervous-system protein OTX1, containing H₁₁, H₅, and an alternating pattern, (HX)₄. The mouse

homeobox protein HOX43 features HA(HP)₄. In the human N-OCT3 (nervous-system specific octamer binding) we observe the pattern HHADH(HP)₂HSHPHQ.

Extensive alternating charge patterns. The two alternating charge runs (+, -)₁₀ = EK(ER)₄EK(ER)₂EKER and (RE)₃ARERD of the disease gene atrophin 1 (DRPLA) are very unusual. A few other examples of similar alternating charge patterns (see refs. 7 and 11) include the large chain of the *Drosophila* developmental transcription factor 5, with (KD)₂(KE)₄; the *Drosophila* female sterile homeotic protein (FSH, see Table 2), with (DR)₄(ER)₃; human HRX protein (see Table 1), with RERD(RE)₂KE; the 101-kDa malaria antigen P-101 (*Plasmodium falciparum*), featuring KE₄(KE)₅(EK)₂E(EK)₃E(EK)₂E(EK)₂(E(EK))₂E₃K₂; the immune system-related RD-protein possessing the precise alternating (+, -)₂₄ sequence [+ is K or R, - is E or D]; and MHC-H2 42 kDa (mouse), which contains an unprecedented (+, -)₂₆ sequence. For possible function and structure implications of hyper-mixed charge runs, see ref. 11.

DISCUSSION

The human triplet repeat disease genes contain multiple long homopeptides, not only of glutamine. Examination of other human, mouse, and particularly *Drosophila* proteins carrying multiple long homopeptides revealed several associated functional and sequence properties. All but a few of the *Drosophila* examples are essential developmental proteins, and a majority of these regulate CNS and/or PNS development. Many of these proteins exhibit a highly anomalous charge distribution with multiple charge clusters interspersed with significantly long uncharged runs and are rich in multiplets and complex repetitive structures. Almost half of the human proteins identified are considered homeotic proteins. Some possible interpretations and implications of these observations are discussed below.

Structural Roles of Long Homopeptides. Why do many developmental regulatory proteins carry multiple long homopeptides, particularly of the small and/or polar amino acids Q, N, H, P, A, S, T, and G? A common proposal suggests that long homopeptides are tolerated nonessential insertions. In some cases, homopeptides of Q, P, S, and T may play a role in transcriptional activation (13–16). In another view, homopeptides may serve as spacer elements between functional domains. As previously indicated, many of these proteins contain two or more charge clusters putatively providing distinct functional domains that interact with DNA, RNA, or proteins and that contribute to intramolecular conformation. The segments linking these domains might often be uncharged polar regions centered on moderate-length homopeptides. In general, if the concentrated charge regions are the essential functional parts of these proteins, the uncharged stretches may play scaffold or hinge roles and thus provide flexibility to the three-dimensional conformation, fine-tuning domain orientation with respect to protein–protein interactions. In this context, excessive expansion of a homopeptide might lead to inappropriate domain organization, resulting in aberrant intramolecular conformation and protein–protein interaction. Moreover, following Perutz *et al.* (8), a drastically extended glutamine (or polar) run in one protein could form a hydrogen bond to a moderate polyglutamine (or polar) peptide and incapacitate the normal counterpart protein in heterozygotes. No complete NMR or crystal structures are yet available for any of the proteins of Tables 1 and 2.

Generation of Multiplets. What is the molecular mechanism that generates high multiplet counts? A high multiplet number indicates an unusual (compared to random) amount of local amino acid iterations. Strand slippage during DNA replication might account for creation of short multiplets and sister-chromatid exchange could be a factor contributing to expan-

sion of long homopeptides. The fact that most proteins of Tables 1 and 2 with significantly large numbers of amino acid multiplets also have significantly high codon multiplet counts argues in favor of strand slippage as a primary mechanism for creating multiplets. A possible explanation for the large proportion of proteins in Tables 1 and 2 that concomitantly have significantly high amino acid and codon multiplet counts is as follows. If long homopeptides usually arise from preexisting short homopeptide/homocodon runs, then those proteins with large numbers of short homopeptides/homocodons would be more likely to acquire one or more long homopeptides. The question still remains as to why the occurrence of large numbers of short and long homopeptides should be so much more frequent in developmental proteins than in other classes of proteins. Is there a difference between essential proteins (probably including most developmental proteins) and nonessential proteins? Are genes expressed at early developmental stages subject to higher levels of mutation? Since most of the long homopeptides are not derived from a single codon it appears that either strand slippage is not the exclusive mechanism for generating amino acid repeats or that many amino acid repeats have been conserved after their creation as homocodon repeats.

Mouse models of neurological disease. A number of mouse models of the HD syndrome have been engineered (18, 19). In all cases where the HD gene is disrupted or turned off, embryonic lethality results in homozygotes, but even in heterozygotes there are behavioral and morphological changes (18), suggesting that the normal protein is critical to development of the animal. In a parallel mouse model on the progression of SCA1 (ataxin 1), polyglutamine expansion leads to rapid neuronal degeneration and death (1). There are also clear asymmetries with respect to HD, SCA-1, and DRPLA diseases in maternal vs. paternal transmission apparently reflected in differences in spermatogenesis, oogenesis, and zygotic development (e.g., refs. 1 and 2).

Proteins with multiple long homopeptides are considerably more abundant in *Drosophila* (about 12% of all available proteins) compared with human and mouse (about 1.6%). Why? On the one hand, there is a well-known positive bias in present collections of *Drosophila* sequences toward developmental proteins. On the other hand, in human and mouse there may be a negative bias in that relatively few nervous system developmental/regulatory genes have been sequenced.

Why do we not observe excessively long homopeptide runs in *Drosophila* or mouse proteins (longest runs are Q₂₅ and Q₂₃, respectively)? Expansions of triplet repeats presumably do occur at some rate in both organisms, but perhaps such mutations are selectively deleterious.

Proposed experimental studies. (i) Although long homopeptides have been deleted from some *Drosophila* proteins [e.g., elimination of the glutamine cluster in Sp1 (14) and in the shaker K⁺ channel protein (20)], the consequences of significantly expanding these repeats have not been experimentally investigated. (ii) Are there diseases (in any tissue) associated with the examples of Table 1? Could expansions of the repeats in the developmental proteins—e.g., engrailed-1, RDC-1, HOX-B3, HOX-B4, OTX-1, and especially the nervous system-specific N-OCT3 (which features Q₂₁)—cause disease? A review of the literature on the human proteins of Table 1 reveals at least two other cases of associated disease. The protein insulin receptor substrate 1 (IRS1) is considered a primary candidate for the site of the defect in insulin action seen in patients with non-insulin-dependent diabetes mellitus (21). Acrosin is the major proteinase present in the acrosome of mature spermatozoa: some abnormalities in acrosin pro-

duce infertility (22). (iii) In the trinucleotide repeat diseases, attention has been focused on expanded polyglutamine tracts. But these proteins contain other long homopeptides, whose expansion might be deleterious. For example, it seems relevant to investigate expansions of the other long homopeptides—e.g., G₂₄ in AR, P₁₁ in the HD gene, S₁₀ in SCA1, etc.—in rodent models and to evaluate phenotypic consequences. Gene sequences from patients suffering infertility or other diseases related to AR could be examined for expansion of the homopeptides G₂₄, P₈, and A₅ as well as Q₂₁. (iv) Most nervous system developmental proteins function widely in many tissues; mutations of these genes apparently induce more severe disorders in the nervous system than in other tissues. Epidemiological surveys and experiments might be designed to search for repeat diseases of other tissues. (v) One or more of the charge clusters could be deleted from the murine homologs of SCA1, DRPLA, HD, or AR proteins and the resulting phenotype characterized.

We are happy to acknowledge valuable discussions concerning triplet repeat diseases with Dr. C. Laird of Seattle and on *Drosophila* development with Drs. M. Krasnow, R. Nusse, and M. Scott of Stanford. Dr. U. Franke was very helpful in reviewing the disease status of the human examples. We also thank Drs. V. Brendel and J. Mrazek for comments on the manuscript. This work was supported in part by National Institutes of Health Grants 2R01GM10452-31 and 5R01HG00335-07 and by National Science Foundation Grant DMS 9403553.

1. La Spada, A. R., Paulson, H. L. & Fischbeck, K. H. (1994) *Ann. Neurol.* **36**, 814–822.
2. Sutherland, G. R. & Richards, R. I. (1995) *Proc. Natl. Acad. Sci. USA* **92**, 3636–3641.
3. Hopkin, K. (1995) *J. NIH Res.* **1**, 45–48.
4. Brendel, V., Bucher, P., Nourbakhsh, I. R., Blaisdell, B. E. & Karlin, S. (1992) *Proc. Natl. Acad. Sci. USA* **89**, 2002–2006.
5. Karlin, S., Ost, F. & Blaisdell, B. E. (1989) in *Mathematical Methods for DNA Sequences*, ed. Waterman, M. (CRC, Boca Raton, FL), pp. 133–157.
6. Green, H. (1993) *Cell* **74**, 955–956.
7. Perutz, M. (1994) *Prot. Sci.* **3**, 1629–1637.
8. Perutz, M. F., Johnson, T., Suzuki, M. & Finch, J. T. (1994) *Proc. Natl. Acad. Sci. USA* **91**, 5355–5358.
9. Scott, K., Blackburn, J. M., Butler, P. J. G. & Perutz, M. (1995) *Proc. Natl. Acad. Sci. USA* **92**, 6509–6513.
10. Karlin, S. (1990) in *Structure & Methods: DNA Protein Complexes & Proteins*, eds. Sarma, R. H. & Sarma, M. H. (Adenine Press, Guilderland, NY), Vol. 2, pp. 171–180.
11. Karlin, S. (1995) *Curr. Opin. Struct. Biol.* **20**, 360–371.
12. Mitchell, P. J. & Tjian, R. (1989) *Science* **245**, 371–378.
13. Brendel, V. & Karlin, S. (1989) *Proc. Natl. Acad. Sci. USA* **86**, 5698–5702.
14. Courey, A. J. & Tjian, R. (1988) *Cell* **55**, 887–898.
15. Su, W., Jackson, S., Tjian, R. & Echols, H. (1991) *Genes Dev.* **5**, 820–826.
16. Gerber, H. P., Seipel, K., Georgiev, O., Höfferer, M., Hug, M., Rusconi, S. & Schaffner, W. (1994) *Science* **263**, 808–811.
17. Brendel, V., Dohlman, J., Blaisdell, B. E. & Karlin, S. (1990) *Proc. Natl. Acad. Sci. USA* **88**, 1536–1540.
18. Nasir, J., Floresco, S. B., O'Kusky, J. R., Diewert, V. M., Richman, J. M., Zeisler, J., Borowski, A., Marth, J. D., Phillips, A. G. & Hayden, M. R. (1995) *Cell* **81**, 811–823.
19. Duyao, M. P., Auerbach, A. B., Ryan, A., Persichetti, F., Barnes, G. T., McNeil, S. M., Ge, P., Vonsattel, J.-P., Gusella, J. F., Joyner, A. L. & MacDonald, M. E. (1995) *Science* **269**, 407–410.
20. Hoshi, T. & Aldrich, R. (1991) *Neuron* **7**, 547–556.
21. Stoffel, M., Espinosa, R., III, Keller, S. R., Lienhard, G. E., Le Beau, M. M. & Bell, G. I. (1993) *Diabetologia* **36**, 335–337.
22. Florke-Gerloff, S., Topfer-Petersen, E., Muller-Esterl, W., Schill, W. B. & Engel, W. (1983) *Hum. Genet.* **65**, 61–67.