

Nucleotide Compositional Constraints on Genomes Generate Alanine-, Glycine-, and Proline-rich Structures in Transcription Factors

Yutaka Nakachi, Toshiyuki Hayakawa, Hiroki Oota, Kenta Sumiyama, Li Wang, and Shintaroh Ueda

Department of Biological Sciences, Graduate School of Science, University of Tokyo

Correlation between amino acid composition and nucleotide composition is examined. Class III POU transcription factors having higher third GC contents showed higher contents of alanine, glycine, and proline residues encoded by GC-rich nucleotides, and vice versa. This correlation was observed even among various types of transcription factors from vertebrates and invertebrates regardless of functional and structural constraints inherent to each protein. Furthermore, reptile class III POU sequences revealed no evolutionary directionality increasing the GC contents from cold- to warm-blooded vertebrates.

Introduction

Proteins have unique amino acid sequences that are specified by their corresponding nucleotide sequences. Since there are 20 amino acids and their corresponding 61 codons, many amino acids are encoded by two or more codons. Such degeneracy of the genetic code permits nucleotide sequences to vary without altering amino acid sequences. As a significant portion of the nucleotide changes at the third codon position are synonymous (silent), the third codon position is considered to directly reflect the degree of nucleotide compositional constraints onto DNAs harboring the sequences.

We have recently discovered a significant correlation between the GC content at the third codon position (the third GC content) and the homopolymeric amino acid repeat ratio of the mammalian class III POU transcription factor genes (Sumiyama et al. 1996): the mammalian Brain-1, Brain-2, and Scip genes have homopolymeric amino acid repeats (sequences without interruptions in the run of a single amino acid residue) including alanine, glycine, and proline, whereas most or all of these repeats are absent from their homologs in nonmammals (amphibians and fish). These characteristic amino acid repeats are well conserved in both position and repeat number among mammals. In contrast, the mammalian Brain-4 gene, like its nonmammalian homolog, has no homopolymeric amino acid repeats. The mammalian Brain-1, Brain-2, and Scip genes containing the homopolymeric amino acid repeats have a higher third GC content. In contrast, the respective nonmammalian homologs lacking the homopolymeric amino acid repeats have a lower third GC content. However, the mammalian Brain-4 gene, like its nonmammalian homolog, has a lower third GC content. There was a clear positive correlation between the homopolymeric amino acid repeat ratio and the third GC content. The amino acids of these characteristic homopolymeric repeats were encoded mainly by codons with a relatively high GC content. These findings indicate that nucleotide

compositional constraints increasing the GC contents (GC pressure) have facilitated the generation of homopolymeric amino acid repeats in mammalian class III POU transcription factors.

To understand how protein sequences have evolved by compositional constraints on genomes, we first determine the class III POU gene structures in reptiles. We then examine correlation between the amino acid composition of proteins and the third GC content using not only the vertebrate class III POU genes, but also various transcription factors from vertebrates and invertebrates.

Materials and Methods

Reptile Class III POU Genes

Recombinant phage clones containing the reptile class III POU genes were isolated from the green anole (*Anolis carolinensis*) genomic library using the POU domain of the chicken Brain-1 gene as a probe. As washing condition in hybridization was lowly stringent (four times for 30 min each in 0.15 M sodium chloride and 0.015 M sodium citrate at 55°C), these clones were identified by hybridization using DNA fragments upstream of the POU domain as probes and partially determining their nucleotide sequences. Nucleotide sequences of the class III POU genes were determined by the dideoxynucleotide chain termination method on both strands. The nucleotide sequence data reported in this paper will appear in the DDBJ/EMBL/GenBank international nucleotide sequence database with the accession numbers AB001868–AB001870.

Sequence Analysis

Transcription factors employed here were as follows; Brain-1 (AB001835, M88299, AB001868, D13045), Brain-2 (L37868, M88300, L27663, AB001869, X64835, Y07905), Brain-4 (X82324, M88301, M84645, AB001870, U17654), and Scip (L26494, X54628, M35205/M72711, X59056, X96422) for the vertebrate class III POU; Evx-1 (D10455, X54239, X60655), Gbx-2 (L39770, L47990), gooseoid (L03395, M81481, M85271, X70471), hairy (L04527, L19314, U36194), HNF-3 β (L10409, L25637), HoxA7 (M17192, M24752), HoxB4 (M26884, M36654), HoxB7 (M16937, X06592), HoxC6 (X12499, X16510), Mash-1 (L11871, M98272, U14587, X53725), Mox2

Key words: GC pressure, POU, transcription factor, reptile.

Address for correspondence and reprints: Shintaroh Ueda, Laboratory of Molecular Biology and Evolution, Department of Biological Sciences, Graduate School of Science, University of Tokyo, 7-3-1 Hongo, Bunkyo-ku, Tokyo 113, Japan. E-mail: suedata@uts2.s.u-tokyo.ac.jp.

Mol. Biol. Evol. 14(10):1042–1049, 1997

© 1997 by the Society for Molecular Biology and Evolution. ISSN: 0737-4038

C

404	ATCCGGGCTTGGAAACCAAGACTCCCTAAAGCCAAGCAAGGGCTCCCATTTGCTCACTGCACACCCCACTTCCACGGATA	484
485	CCGAGGTCCGAGGGCGAGGCTGAAGTCCCTCCCTCTCCCAAGCCAGCCAAATGGGGCCGCACTTCTCTCTGGGACTCCTT	565
566	GGCTTTAACCCCTGCTTTCGCCGGCTGGCTCCCCCTCCCTCGTGGATTTCAATGCTCGGGGGCCGCGGGATTTGGC	646
647	TGCGACGGAGGAGGGAGGGGGAGGGGGAGGGGGGCTCGGCGGGGGCCCTGCCTAGTGCCTACCCCTCTCTGGCTGC	727
728	GGAGAGAGAGATGTCAAAGCGGGCGGGCGGGAGCGCGGCAGAGGAGGCAGCAGCGCGCTCGGCGCGTCCGAGGAGCC	808
809	ATGGCCACGGCGGCTCCAACCCGTACAGCTGCTGGTGCACGCCGAGGCGCGCCGGGGATGCCAGGGCGGGCCCTTC	889
1	M A T A A S N P Y S L L V H A E A A P G M P Q G G P F	27
890	CGCGGCACCCACAGAGTCTCCAGAGCGACTACCTGCAGGGCAACGGGCACCCCTGGGACCCACTGGGTCAACAGC	970
28	R G H H Q K L L Q S D Y L Q G N G H P L G H H W V T S	54
971	CTGAGCGACGCGGGCCCTGGGCTCGAGTTTGGCCGAGCAGCCGACATCAAGCCGGGAGGGAGGACTTGCAGCAGCTG	1051
55	L S D A G P W A S S L A E Q P D I K P G R E D L Q Q L	81
1052	GGGGGGCTTCTGCACCACCGCTCGCCGCTCACCACCACCCAGGCAACGGGGCGTCCGGGGAGGGCGGGACACTTACAA	1132
82	G G L L H H R S P P H H H G N G G V G G G A G H L Q	108
1133	AGCGCTGGTCCAGCAGCCCAACCTCCCGGAATGTTACTCCCAAGCGGGTTCGGCGTGGGAGCCATGCTGGAGCAC	1213
109	S A W S S S P N P P G N V Y S Q G G F G V G A M L E H	135
1214	GGCGGACTACGCCCTCCCGGACCGCCCAATTCCGTCCTCAACAAACAACTGGCCACGGCGCTGCTCCAGAGCCCCAC	1294
136	G G L S P P P T A A N S V P N N N V A T A L L P E P H	162
1295	GACCCTTGAACGCCACCCCGGACCCCTCCGACGAGGAGAGCCCACTCGGACGAGCTGGGACAGTTGCCAAGCAG	1375
163	D P L N S H P G D P S D E E T P T S D E L E Q F A K Q	189
1376	TTCAAGCAGCGGCATCAAGTGGGCTTCAACCAGCCGACGTGGGTTTGGCGCTCGGCACGCTGTACGGCAAGCTTTC	1456
190	F K Q R R I K L G F T Q A D V G L A L G T L Y G N V F	216
	POU-SPECIFIC DOMAIN	
1457	TGCGACAGCAGCTGTCGCGCTTCGAGGCCCTGCAGCTGAGCTTCAAGAACATGTGCAAGCTCAAGCCGCTGCTCAACAAG	1537
217	S Q T T I C R F E A L Q L S F K N M C K L K P L N K	243
1538	TGGCTGGAGGAGGGCCGACTCGTCCACGGGCAGCCACGGGCTCGACAAGATCGCCGCCAGGGCCGGAGAGGAAG	1618
244	W L E E A D S S T G S P T G L D K I A A Q K R K K	270
1619	CGCACCTCCATCGAGTCTCCGTCAAGGGCGTCTGGAGACCCACTTCTCAAGTGCCCAAGCCCGCCGCCAGGAGATC	1699
271	R T S I E V S V K G V L E T H F L K C P K P A A Q E I	297
	POU-HOMEO DOMAIN	
1700	GCCGCCCTCGCCGACAGCTCCAGTGGAGAAGGAAGTCTGTCGCGCTGGTCTGCAACCCGGAGGCAGAAAGGAGAAGCGC	1780
278	A A L A D S L Q L E K E V R V W F C N R R Q K E K R	324
1781	ATGACGTCGCCCGGGGAGAACACGGAGGGGCCCCGCCACGAGGCTACGGGGCGGGGGAGGGGGCCCAAGGCC	1861
325	M T P P G E N N G G A P A H E A Y G A G G G G P K A	351
1862	GACTGCAGGGACCTCTGACCCCGGACTGAGCCAGCCAGCGCCCTGGGGTGGGGGAGACAGAGGCCAGGACCTCCCGTCCG	1942
352	D C R D L *	357
1943	ACTCTTGACTCCGATCTCCAAGACTCCTCTCCGCCAAGCCAACCTCTCTCCAGGCATGGGCACTGCGGCCCTCCC	2023
2024	TCTGGGTGTTTTGGACTTCTCAACAGCCTCAGGCCCTTTGAGCGGAAAGGGAAGGGCCGAGCAGGACATAGCCAGA	2104
2105	ACAAAACCTGGCTATGGCCCTGCTCCTGAAGAGTGGTTCGGTAACGCTAAATAATTTAGGAATCAGGCTCCGTGGATAAACT	2185
2186	TCAGTGGACACTCAAGACTAGGCATCCTCAAGCTGTGGCCCTAACTCCACGTTGTCCAAGGGTCAAGAAATC	2260

FIG. 1 (Continued)

peats in positions 214–229, glycine repeats in positions 31–44, 46–51, and 271–276, and proline repeats in positions 163–168 and 199–204, whereas reptile-specific repeats were glycine repeats in positions 64–69 and 175–179 and glutamine repeats in positions 158–168 and 238–247. Furthermore, there was a wide variation in repeat number of the homopolymeric amino acid repeats common to mammals and reptiles, although all of the amino acid repeats were well conserved among mammals not only in position but also in repeat number. In the Brain-2 gene, the situation was the same as that in the Brain-1 gene. In addition to taxon-specific repeats, there were common repeats including glycine repeats in positions 69–92 and glutamine repeats in positions 139–161 between mammals and reptiles, but a wide variation in repeat number was observed as in the Brain-1 gene (fig. 2B). Surprisingly, the reptile Brain-4 gene also had homopolymeric amino acid repeats, glycine repeats in positions 366–370 (fig. 2C), different from both the mammalian and amphibian homologs with no characteristic amino acid repeats.

Present data on the reptile class III POU genes were consistent with our previous conclusion (Sumiyama et al. 1996): genes having homopolymeric amino acid repeats showed higher third GC contents. The reptile Brain-4, however, did not show a good fit to the plot between the homopolymeric amino acid ratio and the third GC content (see fig. 3A). Moreover, it had another remarkable feature: less amino acid sequence similarity to the mammalian homologs than to that of the amphibian homolog, although mammals are more closely related to reptiles than to amphibians. The phylogenetic tree using the well-conserved POU domain sequence showed the cluster of the mammalian and reptile Brain-4 genes, excluding the possibility that the reptile Brain-4 obtained here was not a homolog to the mammalian and amphibian Brain-4 (data not shown).

Through overall comparison with the other vertebrate homologs, we noticed extraordinary contents of alanine (A), glycine (G), and proline (P) residues in the reptile Brain-4. The codons for these amino acids are GC-rich (GCN, GGN, and CCN, respectively). Con-

A

1 100
 Human MATAASNPLYLPGNSLL--AAGSIVHSDAAGAGGGGGGGGGGGGGGGGGGGGGMPPGSAAVTSGA-----YRGPDSVSKMVSDFM
 Mouse MATAASNPLYQPGNSLL--TAGSIVHSDAAGAGGGGGGGGGGGGGGGGGGGGGMPPGSAAVTSGA-----YRGPDSVSKMVSDFM
 Green anole MATAASNPLYLPGNSLLSAGAIIVHSDAA--AGG-----MQPGSVAVT SVAGGGGGGAGGGGNGNNANNNGYRGPDSVSKMVSDFM
 Zebrafish MATAASNPLYLASSSILSS--GSIVHSDS-----GGGMQGSAAVTSVSGG-----YRGDPT--VKMVSDFM

101 210
 Human -QGAMAASNGGHML SHAHQWVTAL PHAAAAA AAAAAA AV--EASSPWSGS AVGMAGSPQQPPPPPPPPQGP DVKGGAG--RDDLHAGTAL HHRGPPHLGPPPPPHQGHG
 Mouse -QG--AASNGGHML SHAHQWVTAL PHAAAAA AAAAAA AV--EASSPWSGS AVGMAGSPQQPPPPPPPPQGP DVKGGAG--REDLHAGTAL HHRGPPHLGPPPPPHQGHG
 Green anole VQ-----SNGGHML SHAHQWVTAL PHAAAAA AAAAAA AAGSPWSS-----MSGSPQQQQQQQQQQ-----DVKGGGGGREDL-----LHHR--PPHLGPPPP--HQGH--
 Zebrafish -QGAMAASNGGHML SHAHQWVTSL PHAAAAA AAAAAA AVAA--EAGSPWSSSPVGTIGSPQQQ-----DVKNSG--RDDLHSGTAL HNRAP--HLGP-----HQTAYG

211 320
 Human GWGAAAAAAAAAAAAAAAAAHLPSMAGGQQPPQS-----LLYSQPGGFTVNGML SAPPGGGGGGGAGGQAQSLVHPGLVRGDTPELAEH--HHHHHHHAHPHPHPHA
 Mouse GWGAAAAAAAAAAAAAAAAAHLPSMAGGQQPPQS-----LLYSQPGGFTVNGML SAPPGGGGGGGAGGQAQSLVHPGLVRGDTPELAEH--HHHHHHHAHPHPHPHA
 Green anole -WG-----SMAG-QQQQQQQQQQAPLLLYSQPGGFTVNGML SPPGSQALG-----VHPGLVRGDTPELGDHPGHHHHHHHQHHHPHAAH--
 Zebrafish AWG-----STTAHIPS L TGSQQQQ-----FLIYFAPGGFTVNGMHSPP--GS-----QSLVHPGLVRGDTPEL--DHSHHHHHHHHQHHQQAHH--

321 430
POU-SPECIFIC DOMAIN
 Human QGPPHHGGGGGAGPLNSHDPHSDEDTPTSDDLEQFAKQFKQRRRIKLGFTQADVGLALGTYLGNVFSQTTICRFEALQSFKNMCKLKPLL NKWLEEADSSGSPSTSID
 Mouse QGPPHHGGGG--AGPLNSHDPHSDEDTPTSDDLEQFAKQFKQRRRIKLGFTQADVGLALGTYLGNVFSQTTICRFEALQSFKNMCKLKPLL NKWLEEADSSGSPSTSID
 Green anole -----GGGGGGGGGGGLNSHDPHSDEDTPTSDDLEQFAKQFKQRRRIKLGFTQADVGLALGTYLGNVFSQTTICRFEALQSFKNMCKLKPLL NKWLEEADSSGSPSTSID
 Zebrafish -----GVNSHDPHSDEDTPTSDDLEHFAKQFKQRRRIKLGFTQADVGLALGTYLGNVFSQTTICRFEALQSFKNMCKLKPLL NKWLEEADSSGSPSTSID

431 534
POU-HOMEO DOMAIN
 Human KIAAQRKRKRKRTSIEVSVKGALESHFLKCPKPSAQEITSLADSLQLEKEVVRVWF CNRRQKEKRMTPPGIQQTDPDVYSQVGTVSADTPPPH--HGLQTSVQ
 Mouse KIAAQRKRKRKRTSIEVSVKGALESHFLKCPKPSAQEITSLADSLQLEKEVVRVWF CNRRQKEKRMTPPGIQQTDPDVYSQVGTVSADTPPPH--HGLQTSVQ
 Green anole KIAAQRKRKRKRTSIEVSVKGALESHFLKCPKPSAQEITSLADSLQLEKEVVRVWF CNRRQKEKRMTPPGIQQAADDVYSQVGNVSADTPPPH--HGLQGGVQ
 Zebrafish KIAAQRKRKRKRTSIEVSVKGALESHFLKCPKPSAQEITSLADSLQLEKEVVRVWF CNRRQKEKRMTPPGVPQ--TPEDVYSQVGNVSADTPPPSMDCKRMFSET

B

1 60
 Human MATAASNHYSLTSSASIVHAEPGAMQAGGYREAQSL--VQGDYALQSNHGPHL SHAH
 Mouse MATAASNHYSLTSSASIVHAEPGGMQAGGYREAQSL--VQGDYALQSNHGPHL SHAH
 Rat MATAASNHYSLTSSASIVHAEPGGMQAGGYREAQSLVQGDYALQSNHGPHL SHAH
 Green anole MATTASNHYSLLAASSPMVHAEPGSMQPGAG--YRDA-----VQADYALQSNHGPHL SHAH
 Xenopus MATTASNHYNLGSGSSIVHADP--GGMQQAQ--YRDAQTL--VQSDYT--LQSNHGPHL SHAH
 Zebrafish MATTASNHYNILTSSPSIVHSEP--GSMQATA--YRDAQTL--LQSDYS--LQSNHGPHL SHAH

61 170
 Human QWITALSHGGGGGGGGGGGGGGGGGGGG-----DGSPWSTSP LQPDIKPSVVVQGGGRDELHGPALQOQ--H--QOQQQQQQQQQQQQQQQQQQQ--RPPHLVHA
 Mouse QWITALSHGGGGGGGGGGGGGGGGGGGG-----DGSPWSTSP LQPDIKPSVVVQGGGRDELHGPALQOQ--H--QOQQQQQQQQQQQQQQQQQQQ--RPPHLVHA
 Rat -WITALSHGGGGGGGGGGGGGGGGGGGG-----DGSPWSTSP LQPDIKPSVVVQGGGRDELHGPALQOQ--H--QOQQQQQQQQQQQQQQQQQQQ--RPPHLVHA
 Green anole QWITAAALSHGGGGGGGGGGGGGGGGGGGGSSGD--SPWSTSP---DIPKPSV--QAGGRDDL-----QOQQHHQQQQQQQQ-----GRPPHLVHA
 Xenopus QWITALSHG-----DGAPWATSP LQOQDIKPTV--Q--SSR--DELHVSGLTQ---H--Q-----SRAPHLVHA
 Zebrafish QWITALSHG-----EGPWSSSPLGEQDIKPAV--Q--SPR--DEMHNSSNLQ---H--Q-----SRPPHLVHT

171 280
 Human ANHH--PGGA--WR TAAAAAHL PPSMGASNGG---LLYSQP--SFTVNGMLGAG--GQPAGLHHHGLRDAHDEP-----HHADHHPHPH--SHPHQQ-----PP
 Mouse ANHH--PGGA--WR SAAAAAHL PPSMGASNGG---LLYSQP--SFTVNGMLGAG--GQPAGLHHHGLRDAHDEP-----HHADHHPHPH--SHPHQQ-----PP
 Rat ANHH--PGGA--WR SAAAAAHL PPSMGASNGG---LLYSQP--SFTVNGMLGAG--GQPAGLHHHGLRDAHDEP-----HHADHHPHPH--SHPHQQ-----PP
 Green anole GSHAAVAAA VAWR--TGGS AHL PPGMAAANGGAQQGGL LYSQPPPGFTVNGML--GSGQP--GMHHHGLREAHEPPPPPPPPPPHPHDHL-----SQOQQQQQQHAPP
 Xenopus HGNNH--GPGA--WRSTGST--HLS--SMASNG--QG--LLYSQP--SFTVNGMINPGSQ--GIHHHGLRDSHDD-----HHGDHG--HQQVSAQQQHS LQLQ----
 Zebrafish HGNNH--DSRA--WR--TTAAHIP--SMATSNG--QS--LIYSQP--SFSVNGLI--PGSQG--GIHHMSRDAHED-----HHSPLSDHG--HPP--SQ--HQ--HQSQS---

281 390
POU-SPECIFIC DOMAIN
 Human P P P P P Q G P - P G H P G A H D P H S D E D T P T S D D L E Q F A K Q F K Q R R I K L G F T Q A D V G L A L G T Y L G N V F S Q T T I C R F E A L Q L S F K N M C K L K P L L N K W L E E A D S S G S P T S I D K I A
 Mouse P P P P P Q G P - P G H P G A H D P H S D E D T P T S D D L E Q F A K Q F K Q R R I K L G F T Q A D V G L A L G T Y L G N V F S Q T T I C R F E A L Q L S F K N M C K L K P L L N K W L E E A D S S G S P T S I D K I A
 Rat P P P P P Q G P - P G H P G A H D P H S D E D T P T S D D L E Q F A K Q F K Q R R I K L G F T Q A D V G L A L G T Y L G N V F S Q T T I C R F E A L Q L S F K N M C K L K P L L N K W L E E A D S S G S P T S I D K I A
 Green anole P P H H H H P H P A H H P H H E A H S D E D T P T S D D L E Q F A K Q F K Q R R I K L G F T Q A D V G L A L G T Y L G N V F S Q T T I C R F E A L Q L S F K N M C K L K P L L N K W L E E A D S S G S P T S I D K I A
 Xenopus -----GGHQD--HSDEDTPTSDDLEQFAKQFKQRRRIKLGFTQADVGLALGTYLGNVFSQTTICRFEALQSFKNMCKLKPLL NKWLEEADSSGSPSTSIDKIA
 Zebrafish -----HHD--HSDEDTPTSDDLEQFAKQFKQRRRIKLGFTQADVGLALGTYLGNVFSQTTICRFEALQSFKNMCKLKPLL NKWLEEADSSGSPSTSLDKIA

391 486
POU-HOMEO DOMAIN
 Human A Q G R K R K R T S I E V S V K G A L E S H F L K C P K P S A Q E I T S L A D S L Q L E K E V V R V W F C N R R Q K E K R M T P P G G T L P G A E D V Y G S R D T P P - H H G V Q T P V Q
 Mouse A Q G R K R K R T S I E V S V K G A L E S H F L K C P K P S A Q E I T S L A D S L Q L E K E V V R V W F C N R R Q K E K R M T P P G G T L P G A E D V Y G S R D T P P - H H G V Q T P V Q
 Rat S Q G R K R K R T S I E V S V K G A L E S H F L K C P K P S A Q E I T S L A D S L Q L E K E V V R V W F C N R R Q K E K R M T P P G G T L P G A E D V Y G S R D T P P - H H G V Q T P V Q
 Green anole A Q G R K R K R T S I E V S V K G A L E S H F L K C P K P S A Q E I T S L A D S L Q L E K E V V R V W F C N R R Q K E K R M T P P G G T L P G A E D V Y G A S R D T P P H H G V Q T P V Q
 Xenopus A Q G R K R K R T S I E V S V K G A L E S H F L K C P K P S A Q E I T S L A D S L Q L E K E V V R V W F C N R R Q K E K R M T P P G G T L P G A E D V Y G A S R D T P P - H L G V Q T S V Q
 Zebrafish A Q G R K R K R T S I E V S V K G A L E S H F L K C P K P A S E I T S L A D S L Q L E K E V V R V W F C N R R Q K E K R M T P P G G L P G T E D V Y G ---DTPP--HHGVQTPVQ-

FIG. 2.—Sequence alignment within the homologs of the class III POU genes. A, Brain-1. B, Brain-2. C, Brain-4. The POU-specific and POU-homeo domains are indicated with shadowed boxes. Gaps are represented by hyphens.

C

Human		1	MATAASNYPYLSSTSLVHADSA-GMQQGSFPR-NPQKLLQSDYLQGVPSNGHPLGHHWV	60
Mouse			MATAASNYPYLSSSSLVHADSA-GMQQGSFPR-NPQKLLQSDYLQGVPSNGHPLGHHWV	
Rat			MATAASNYPYLSSSSLVHADSA-GMQQGSFPR-NPQKLLQSDYLQGVPSNGHPLGHHWV	
Green anole			MATAASNYPYSL-----LVHAEAAPMPQGGPFRGHQKLLQSDYLQGG---NGHPLGHHWV	
Xenopus			MATAASNYPYLSSSSLVHADSA-VMQQGSFPR-NPQKLLQSDYLQGVPCNGHPLGHHWV	
Human	61		TSLSDGGPWSSTLATSPLDQQDVKPGREDLQ-LGAIHHRSPhVAHHS-----HTNHPNAGWASPAPNSITSSGQPLNVYSQPGFTVSGMLEHGGLTPPPAAASAQS	170
Mouse			TSLSDGGPWSSTLATSPLDQQDVKPGREDLQ-LGAIHHRSPhVAHHS-----HTNHPNAGWASPAPNSITSSGQPLNVYSQPGFTVSGMLEHGGLTPPPAAASTQS	
Rat			TSLSDGGPWSSTLATSPLDQQDVKPGREDLQ-LGAIHHRSPhVAHHS-----HTNHPNAGWASPAPNSITSSGQPLNVYSQPGFTVSGMLEHGGLTPPPAAASTQS	
Green anole			TSLSDAGPWASSLA-----EQDPDKPGREDLQQLGGLLHHRSPPHHHHNGGVGGAGHLQSAWSSSPNP-----PGNVYSQGGFVGAMLEHGGLSPPTAANSVP	
Xenopus			TSLSDANPWSSSLASSPLDQQDVKPGREDLQ-LGAIHHRSPhVNHHS-----HTNHPNAGWASPAPNSITSSGQPLNVYSQPGFTVSGMLDHGELTPPLPAGTTQS	
POU-SPECIFIC DOMAIN				
Human	171		LHPVLR-----PPDHGELGSHHCQDHSDEETPTSDLEQFAKQFKQRRIKLGFTQADVGLALGTLYGNVFSQTTICRFEALQLSFKNMCKLKPLLKNKWLLEADSTGSPSTSI	280
Mouse			LHPVLR-----PPDHGELGSHHCQDHSDEETPTSDLEQFAKQFKQRRIKLGFTQADVGLALGTLYGNVFSQTTICRFEALQLSFKNMCKLKPLLKNKWLLEADSTGSPSTSI	
Rat			LHPVLR-----PPDHGELGSHHCQDHSDEETPTSDLEQFAKQFKQRRIKLGFTQADVGLALGTLYGNVFSQTTICRFEALQLSFKNMCKLKPLLKNKWLLEADSTGSPSTSI	
Green anole			NNNVATALLPEPHDPLNSHP-GDPSDEETPTSDLEQFAKQFKQRRIKLGFTQADVGLALGTLYGNVFSQTTICRFEALQLSFKNMCKLKPLLKNKWLLEADSTGSPSTGL	
Xenopus			LHPVLR-----PNDHVDLGSHHHCQDHSDEETPTSDLEQFAKQFKQRRIKLGFTQADVGLALGTLYGNVFSQTTICRFEALQLSFKNMCKLKPLLKNKWLLEADSTGNPTSI	
POU-HOMEO DOMAIN				
Human	281		DKIAAQRKRKRRTSIEYSVKGVLETHFLKCPKPAAQEISSLADSLQLEKEVVRVWFNRRQKEKRMTPGDDQ----PHEVYSHTVKTD--TSCHDL	378
Mouse			DKIAAQRKRKRRTSIEYSVKGVLETHFLKCPKPAAQEISSLADSLQLEKEVVRVWFNRRQKEKRMTPGDDQ----PHEVYSHTVKTD--ASCHDL	
Rat			DKIAAQRKRKRRTSIEYSVKGVLETHFLKCPKPAAQEISSLADSLQLEKEVVRVWFNRRQKEKRMTPGDDQ----PHEVYSHTVKTD--ASCHDL	
Green anole			DKIAAQRKRKRRTSIEYSVKGVLETHFLKCPKPAAQEISSLADSLQLEKEVVRVWFNRRQKEKRMTPGGENNGAPAEHAYGAGGGGGPKADCRDL	
Xenopus			DKIAAQRKRKRRTSIEYSVKGVLETHFLKCPKPAALETSLADSLQLEKEVVRVWFNRRQKEKRMTPGDDQ----QHEVYSHTVKTD--TSCNEL	

FIG. 2 (Continued)

versely, the number of arginine residues, the remaining amino acid encoded by GC-rich codons, did not vary among the vertebrate Brain-4 genes (14 residues for all the vertebrate Brain-4 genes). The Brain-2 gene also showed no variation in number of arginine residues (16 for all the vertebrates). We found some variations in arginine number in both Brain-1 and Scip genes, but their difference is only one residue (14 and 15 for human/mouse/green anole and zebrafish Brain-1, respectively, and 15 and 16 for human/mouse/rat/*Xenopus* and zebrafish Scip, respectively). We thus used the total content of A, G, and P residues (AGP content) to investigate a correlation between the third GC content and the amino acid composition among the vertebrate class III POU genes. As the third GC content of the entire region is shown to be nearly equal to that of the POU domain with no homopolymeric amino acid repeats (Sumiyama et al. 1996), we used the third GC content of the entire region so as to make comparisons among a wide variety of transcription factors possible, as mentioned below. We found a significant correlation between the third GC and AGP contents (correlation coefficient was 0.82), as shown in figure 3B.

Phylogenetic analysis indicates that Brain-1, Brain-2, Brain-4, and Scip genes already existed in the common ancestor of vertebrates (Sumiyama et al. 1996). These genes are not tandemly located on the same chromosome, at least in mammals; the mammalian Brain-4 gene is located on the evolutionarily well-conserved X chromosome, while the other class III POU genes are on autosomes (Avraham et al. 1993; Xia et al. 1993; Atanasoski et al. 1995; de Kok et al. 1995). The present data thus suggest that (1) the ancestral class III POU gene possesses no enrichment of A, G, or P residues, (2) enrichment of AGP residues including generation of the homopolymeric amino acid repeats occurred inde-

pendently both in particular lineages and in particular class III POU genes (AGP enrichment in the Brain-1 and Brain-2 genes occurred both in the common ancestor of amniotes and independently in each lineage of amniotes, whereas that in the Brain-4 gene occurred only in reptiles). Moreover, we conclude that the hypothesis based on evolutionary directionality increasing the GC contents from cold- to warm-blooded vertebrates is incorrectly drawn due to sparse sequence data on cold-blooded vertebrates except for amphibians. In fact, there are only five reptile genes used in figure 3 of Bernardi and Bernardi (1991), and the GC contents of those appear to be nearly equal to those of the mammalian homologs.

A similar situation holds for other transcription factors of vertebrates: a wide variation of AGP content and a nearly equal arginine content. We found a clear tendency for the AGP content to increase relative to the third GC content in each transcription factor (data not shown). We thus plotted the AGP content against the respective third GC content for all the vertebrate transcription factors studied (fig. 4A). To our surprise, there was a clear positive correlation (correlation coefficient was 0.72) regardless of functional and structural constraints inherent in each protein. We also analyzed data from both vertebrates and invertebrates (fig. 4B) and again found a positive correlation (correlation coefficient was 0.71). No arginine-rich region was found in any of the transcription factors examined. This observation is compatible with the frequency distribution of homopolymeric amino acid repeats of proteins in general (Green and Wang 1994).

Present results provide a general picture for protein structure and its evolution: amino acid compositions are under profound influence of nucleotide compositional constraints on genome DNAs harboring coding sequen-

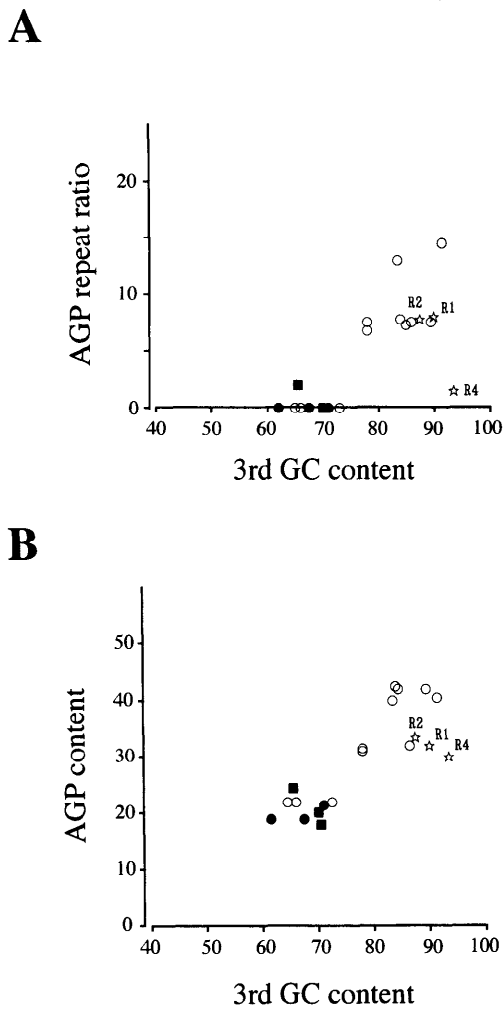


FIG. 3.—*A*, Plot of the homopolymeric alanine/glycine/proline amino acid repeat ratio (AGP repeat ratio) of the vertebrate class III POU transcription factors against the respective third GC content. *B*, Plot of the alanine/glycine/proline amino acid content (AGP content) of the vertebrate class III POU transcription factors against the respective third GC content. ○: mammals, ☆: reptiles, ●: amphibians, ■: fish. R1, R2, and R4 represent the green anole Brain-1, Brain-2, and Brain-4 genes, respectively. Homopolymeric amino acid repeats are defined as sequences consisting of more than four consecutive identical amino acid residues without interruptions.

es. As a result, the ratio of A, G, and P residues linearly correlates with the degree of nucleotide compositional constraints increasing the GC contents, and changes in nucleotide compositional constraints have caused concomitant alterations in amino acid compositions through evolution.

A-, G-, and P-rich sequences are identified as transcriptional activation domains of transcription factors (Mermod et al. 1989; Mitchell and Tjian 1989; Licht et al. 1990; Tanaka, Clouston, and Herr 1994; Catron et al. 1995). In fact, a transcription factor artificially fused with homopolymeric proline repeats significantly modulates its transcriptional activation (Gerber et al. 1994). We therefore suggest that enrichment of A, G, and P residues in transcription factors caused by GC pressure should have a profound influence on diversification of

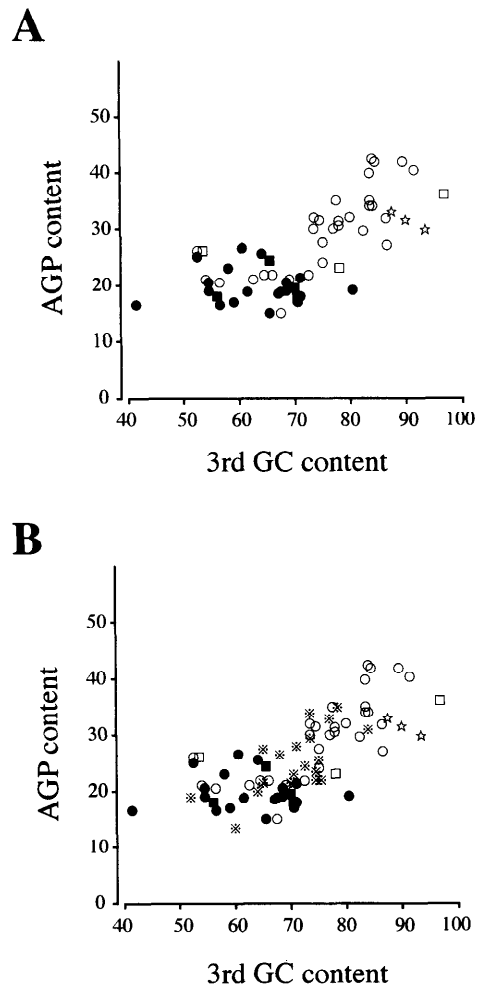


FIG. 4.—Correlation between the alanine/glycine/proline amino acid content (AGP content) of transcription factors and the respective third GC content. *A*, Vertebrates. *B*, Vertebrates plus invertebrates. ○: mammals, □: avians, ☆: reptiles, ●: amphibians, ■: fish, ※: *Drosophila*. Correlation coefficient for invertebrates alone was 0.67.

gene regulation mechanisms (fig. 5). A possible functional difference in transcription factors caused by nucleotide compositional constraints will be the subject of forthcoming studies.

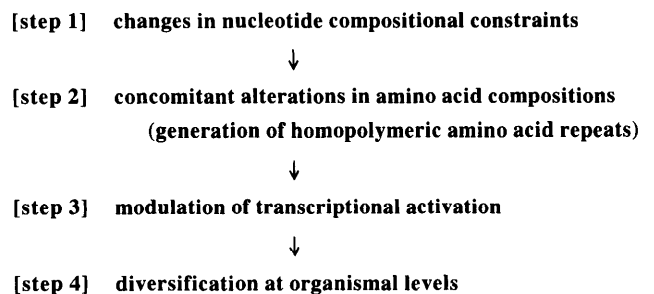


FIG. 5.—Proposed hypothesis on molecular mechanisms for diversification at the organismal level caused by genomic evolution. Note that this hypothesis is applicable not only to species diversification caused by functional differentiation among orthologous genes, but also to functional differentiation among paralogous (duplicated) genes.

Acknowledgments

We thank S. Yokoyama for the reptile library. We also gratefully acknowledge the thoughtful criticism of N. Saitou and the anonymous reviewers. This study was supported by grants from the Ministry of Education, Science, Sports, and Culture of Japan.

LITERATURE CITED

- ATANASOSKI, S., S. S. TOLDO, U. MALIPIERO, E. SCHREIBER, R. FRIES, and A. FONTANA. 1995. Isolation of the human genomic Brain-2/N-Oct 3 gene (*POUF3*) and assignment to chromosome 6q16. *Genomics* **26**:272–280.
- AVRAHAM, K. B., B. C. CHO, D. GILBERT et al. (15 co-authors). 1993. Murine chromosomal location of four class III POU transcription factors. *Genomics* **18**:131–133.
- BERNARDI, G. 1995. The human genome: organization and evolutionary history. *Annu. Rev. Genet.* **29**:445–476.
- BERNARDI, G., and G. BERNARDI. 1991. Compositional properties of nuclear genes from cold-blooded vertebrates. *J. Mol. Evol.* **33**:57–67.
- BERNARDI, G., B. OLOFSSON, J. FILIPSKI, M. ZERIAL, J. SALINAS, G. CUNY, M. MEUNIER-ROTIVAL, and F. RODIER. 1985. The mosaic genome of warm-blooded vertebrates. *Science* **228**:953–958.
- CATRON, K. M., H. ZHANG, S. C. MARSHALL, J. A. INOSTROZA, J. M. WILSON, and C. ABATE. 1995. Transcriptional repression by *Msx-1* does not require homeodomain DNA-binding sites. *Mol. Cell. Biol.* **15**:861–871.
- DE KOK, Y. J. M., S. M. VAN DER MAAREL, M. BITNER-GLINDZICZ, I. HUBER, A. P. MONACO, S. MALCOLM, M. E. PEMBREY, H. H. ROPERS, and F. P. M. CREMERS. 1995. Association between X-linked mixed deafness and mutations in the POU domain gene *POU3F4*. *Science* **267**:685–688.
- GERBER, H., K. SEIPEL, O. GEORGIEV, M. HÖFFERER, M. HUG, S. RUSCONI, and W. SCHAFFNER. 1994. Transcriptional activation modulated by homopolymeric glutamine and proline stretches. *Science* **263**:808–811.
- GREEN, H., and N. WANG. 1994. Codon reiteration and the evolution of proteins. *Proc. Natl. Acad. Sci. USA* **91**:4298–4302.
- IKEMURA, T., and S. AOTA. 1988. Global variation in G+C content along vertebrate genome DNA. *J. Mol. Biol.* **203**:1–13.
- KADI, F., D. MOUCHIROUD, G. SABEUR, and G. BERNARDI. 1993. The compositional patterns of the avian genomes and their evolutionary implications. *J. Mol. Evol.* **37**:544–551.
- LICHT, J. D., M. J. GROSSEL, J. FIGGE, and U. M. HANSEN. 1990. Drosophila Krüppel protein is a transcriptional repressor. *Nature* **346**:76–79.
- MERMOD, N., E. A. O'NEILL, T. J. KELLY, and R. TJIAN. 1989. The proline-rich transcriptional activator of CTF/NF-1 is distinct from the replication and DNA binding domain. *Cell* **58**:741–753.
- MITCHELL, P. J., and R. TJIAN. 1989. Transcriptional regulation in mammalian cells by sequence-specific DNA binding proteins. *Science* **245**:371–378.
- SUMIYAMA, K., K. WASHIO-WATANABE, N. SAITOU, T. HAYAKAWA, and S. UEDA. 1996. Class III POU genes: generation of homopolymeric amino acid repeats under GC pressure in mammals. *J. Mol. Evol.* **43**:170–178.
- TANAKA, M., W. M. CLOUSTON, and W. HERR. 1994. The Oct-2 glutamine-rich and proline-rich activation domains can synergize with each other or duplicates of themselves to activate transcription. *Mol. Cell. Biol.* **14**:6046–6055.
- THOMPSON, J. D., D. G. HIGGINS, and T. J. GIBSON. 1994. CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice. *Nucleic Acids Res.* **22**:4673–4680.
- XIA, Y. R., B. ANDERSON, M. MEHRABIAN, A. T. DIEP, C. H. WARDEN, T. MOHANDAS, R. J. MCEVILLY, M. G. ROSENFELD, and A. J. LUSIS. 1993. Chromosomal organization of mammalian POU domain factors. *Genomics* **18**:126–130.

NARUYA SAITOU, reviewing editor

Accepted June 30, 1997