Origin and Evolution of Influenza Virus Hemagglutinin Genes

Yoshiyuki Suzuki and Masatoshi Nei

Institute of Molecular Evolutionary Genetics and Department of Biology, The Pennsylvania State University

Influenza A, B, and C viruses are the etiological agents of influenza. Hemagglutinin (HA) is the major envelope glycoprotein of influenza A and B viruses, and hemagglutinin-esterase (HE) in influenza C viruses is a protein homologous to HA. Because influenza A virus pandemics in humans appear to occur when new subtypes of HA genes are introduced from aquatic birds that are known to be the natural reservoir of the viruses, an understanding of the origin and evolution of HA genes is of particular importance. We therefore conducted a phylogenetic analysis of HA and HE genes and showed that the influenza A virus HA genes. The rate of amino acid substitution for A virus HAs from duck, a natural reservoir, was estimated to be 3.19×10^{-4} per site per year, which was slower than that for human and swine A virus HAs but similar to that for influenza B and C virus HAs (HEs). Using this substitution rate from the duck, we estimated that the divergences between different subtypes of A virus HA gene diverged from several thousand to several hundred years ago. In particular, the earliest divergence time was estimated to be about 2,000 years ago. Also, the A virus HA gene diverged from the B virus HA gene about 4,000 years ago and from the C virus HE gene about 8,000 years ago. These time estimates are much earlier than the previous ones.

Introduction

Influenza viruses are members of the viral family Orthomyxoviridae and have a segmented, single-stranded, and negative-sense RNA genome in an enveloped virion (Smith, Andrewes, and Laidlaw 1933). The genome encodes envelope glycoproteins, matrix proteins, nonstructural proteins, nucleoproteins, and polymerase proteins. According to the antigenic properties of matrix proteins or nucleoproteins, influenza viruses are classified into types A, B, and C. Influenza A viruses cause epidemics and pandemics of influenza in mammals and birds, and aquatic birds are known to be the natural reservoir of these viruses (Slemons et al. 1974; Webster et al. 1978; Hinshaw, Webster, and Turner 1980). Influenza B and C viruses are isolated mainly from humans and are less pathogenic than influenza A viruses. Phylogenetic analyses of nucleoproteins and polymerase proteins have indicated that influenza A and B viruses are more closely related to each other than to influenza C viruses (Gammelin et al. 1990; Krossoy et al. 1999; Cox et al. 2000).

Hemagglutinin (HA) is the major envelope glycoprotein of A and B viruses, and hemagglutinin-esterase (HE) in C viruses is a protein homologous to HA. HA (HE) is cleaved into the signal peptide (about 20 amino acids in influenza A viruses), protein HA1 (HE1) (about 320 amino acids), and protein HA2 (HE2) (about 220 amino acids) when mature proteins are produced (fig. 1). HA1 (HE1) is a receptor-binding protein and the major target of immune responses, whereas HA2 (HE2) is an anchor protein of the envelope and mediates fusion of the envelope and the cellular endosomal membrane. Influenza A virus HA genes are classified into 15 subtypes (H1–H15), according to their antigenic properties

Key words: influenza virus, hemagglutinin, hemagglutinin-esterase, rate of amino acid substitution, divergence time.

Mol. Biol. Evol. 19(4):501-509. 2002

© 2002 by the Society for Molecular Biology and Evolution. ISSN: 0737-4038

(WHO Memorandum 1980), whereas B and C virus HA (HE) genes are not classified into subtypes. Because influenza A virus pandemics in humans appear to occur when new subtypes of HA genes are introduced from aquatic birds, an understanding of the origin and evolution of HA genes is of particular importance.

From the phylogenetic analyses of A and B virus HA genes, Webster et al. (1992) suggested that the divergence between these two genes occurred later than the divergences between some subtypes of A virus HA genes. However, the reliability of their conclusion is unclear because they did not use any outgroup to root their phylogenetic tree. Saitou and Nei (1986) estimated the earliest divergence time between subtypes of A virus HA genes to be 200–300 years ago, using the rate of amino acid substitution for human A virus HAs. This estimate is also unreliable because the natural reservoir of these viruses is aquatic birds, and the evolutionary rate is known to be slower in birds than in humans (Saitou and Nei 1986; Bean et al. 1992; Schafer et al. 1993; Makarova et al. 1999; Reid et al. 1999; Suarez 2000).

The purpose of this paper is to study the evolutionary relationships of influenza A, B, and C virus HA (HE) genes. We are also interested in estimating the divergence times between these genes.

Materials and Methods

Phylogenetic Analyses

For constructing a phylogenetic tree for influenza A, B, and C virus HA (HE) genes, we used amino acid sequences because they are known to give more reliable results than nucleotide sequences when the sequence divergence is high (Nei and Kumar 2000, pp. 17–32). We also used only the HA2 (HE2) region of HA (HE) because in the multiple alignment for the entire region of A, B, and C virus HAs (HEs), the signal peptide and HA1 (HE1) protein regions generated many gaps and only the HA2 (HE2) region appeared reliable (fig. 1) (Nakada et al. 1984).

Amino acid sequences of influenza A, B, and C virus HA2s (HE2s) were collected from the international

Address for correspondence and reprints: Yoshiyuki Suzuki, Institute of Molecular Evolutionary Genetics, The Pennsylvania State University, 328 Mueller Laboratory, University Park, Pennsylvania 16802. E-mail: yis1@psu.edu.

18 SSVSSFEH SSVTHFER SSVTHFER SSTTHFER SSTTHFER SSTTFFER SSTRFFER SSACSTTR	5 5 5 6 6 6 6 6 7 7 7 7 7 7 7 7 7 7 7 7 7	KEPNKLER KEPNKLER KEPNKLER KEPNKLER KEPNKLER KEPNKLER HEPSELER HEPSELER HEPSELER HEPSELER KEPFQUER GENQUER GENQUER GAMELHI	
EELREQL SELREQL ASELRELU ASLRSLU ASLRSLU ASLRSLU ASLRSLU SELRELL GELRHLF GELRHLF GELRHLF GELRHLF GELRHLF GELRHLF GELRHLF GELRHLF GELRHLF GELRHLF SVITYELA	YYWTLLK FFSWTLD FFSWTLD FFSWTLD FFWTLUF FFWTLJLL FFWTLJLL FFWTLJLL FFWTLJLL FFWTLJLL FFWTLJLL FFWTLJLL FFWTLJLL FFWTLJLL FFWTLJL FFWTJL FFWTLJ	IIQFTAVG ITQFGAVG ITQFGAVG ITQFGAVG ITQFBAVG ITQFBAVG ITQFBAVG ITQFBAVG ITQFBAVD ITQFBAVD ITQFBAVD ITQFBAVD ITQFBAVD ITQFBAVD ITQFBAVD ITQFBAVD ITQFBAVD ITQFBAVD ITQFBAVD ITQFBAVD ITTFFB	
DF IDY PARTY Control of the second s	0.00.00.00.00.00.00.00.00.00.00.00.00.0	NEVIEKA SUVIEK	
LCYPG		NGTTNKVI DGTTNKVI DGTTNGKLI DGTDGKLI DGTDGKLI DGTTNKVI DGTTTNKVI DGTTTNKVI DGTTTNKVI DUNONKII DVITTSKI DUNONKII DVITTSKI	
TPINSENG KENTRUCK RSKAYS- RSKAYS- KASPAND KRFADNG RESAPDG RESAPDG RESAPDG RESAPDG RESAPDG KFNPTNG KFNPTNG KRNSSD- CRPAND- CKPATS- VKPATS- WSPHAAI		KSTONAL ESTORAL KSTOALL KSTOAL KSTOAL KSTOSAL KSTOBAL KSTOAL CSTORAL CSTORAL CSTORAL CSTORAL CSTORAL CSTORAL	688 688 688 688 688 688 688 6681610000000000
RSMSYIVE REWSYIRE REWSYIRE REWSYIRE REWSYIRE REWSYIRE REMSYIRE REMSYIRE REMSYIRE REMSYIRE REMSYIRE REMSYIRE REMSYIRE REMSYIRE REMSYIRE REMSYIRE		GSGYAAD GSGYAAD GSGYAAD GSGYAAD GSGYAAD GSGYAAD GSGYAAD GSGYAAD GSGAAD GSGAAD GSGAAD GSGTAAD GSGTAAD GSGTAAD GSGTAAD GSGTAAD GSGTAAD	MCSNGSL/ MCSNGSL/ ACCKGNTI ACCKGNTI ACCKGNU MCSNGSL/ MCSNGSL/ MCSNGSL/ MCSNGSL/ MCSNGSL/ MCSNGSL/ MCSNGSL/ CIKNGNM ACSNGSL/ GCONGNT/ GCONGNT/ CVKNGNM ACSNGSL/ ACCNGNT/ CVKNGNM ACSNGSL/ ACCNGNT/ CVKNGNM
PU		YHHONBO PRHONBO PRHONBO YHHSNDO PRHONBO FRHONBO FRHONBO FRHONBO FRHONBO FRHONBO FRHONBO FRHONBO FRHONBO FRHONBO FRHONA	LIGAISFW LIGAISFW LLAFILW LLAFLAFLW LLAFLAFLA LAJAULAFLAFLA LAJAULAFLAFLA LAJAULFF LIGGFFF LIGGFFF LIGGFLFW LALEVLW LALEVLW LALEVLW LALEVLW
UPEC DPLI UPEC DPLI UPEC DDFL UPEC DDFL UPEC DPLI UPEC D	ABR P KVRL ABR P KVNC SSR P WVRC SSR P WVRC SSR P WVRC SSR P VVRC ABR P 2 VVRC SSR P L VVRC SSR P L VVRC SSR P KVNC SSR P KVNC	* MUDGWYG- MVDGWYG- MVDGWYG- MUDGWYG- LIDGWYG- LIDGWYG- LUDGWYG- LUNGWYG- LIDGWYG- L	SSLVLLVV SSLVLLVV SSFLLCTVN SSLLLAN SSLLLLAN SSLVLLAN ASLCLANIN ASLCLANIN ASLCLANIN ASLCLANIN ASLCLANIN ASLCLANIN ASLCLANIN ASLCLANIN ASLCLANIN ASLCLANIN ASLCLANIN ASLCANINA SSLVVLLAN SSLAVLLAN SSLAVLLAN LGLATTAN
IAGMLLG IAGMLLG IIDALLG IIDALLG IIDALLG IIDALLG IIDALLG IIDALLG IIDALLG IIDALLG IILOTLTG ILDALLG ILDALLG IILOTLTG IILOTLTG IILOTLTG IILOTALG IILOTA	RRFTPELI RRFTPELI CTSVTPHIL CTSVTPHIL CTSVTPHIL CTSVTPELI CTSVTPELI PAKEPELI RTFTPELI RSFKPHIL RSFKPHIL RSFKPHIL RSFKPHIL RSFKPHIL RSFKPHIL COSTSPERI SPTGHIS	IEGGWTCR IEGGWTCR IERGGWCG IERGGWCG IERGGWCG IERGGWCG IEGGWPCG IEG	AIYSTVA LAIYSTVA LAIYSTVA LAISFAI LUNSFSI LUNSFSA LINFSFGA LINFSFGA LINFSFGA LINFSFGA LINFSFGA LINFSFGA LINFSFGA LINFSFGA LINFSFGA
LQLGRCN LLGLGRNCT LLGLARNCT LLGLARNCT VJDQTCD VJDQTCD VJLDDCS VLLDDCS VLLDDCS VLLDDCS VVDGQCCH VVDQCCH VVDQCH VVD	5VVTSNYN VVSTKRSL VVSTKRSL VVSTKRSL VVSTKRSL VVSTKRSL VVSTKRSL VVSTKSL VVSTKSL VVSTKSL VVSTKSL VVSTKSL VVSTKSL VVSTKSL VS	IA2 FEGIAGF FEGIAGF FEGIAGF FEGIAGF FEGIAGF FEGIAGF FEGIAGF FEGIAGF FEGIAGF FEGIAGF FEGIAGF FEGIAGF FEGIAGF FEGIAGF FEGIAGF	11071-0 11071-0 11071-0 11071-0 11071-0 11071-0 11177-0 111777-0 111777-0 111777-0 111777-0 111777-0 111777-0 111777-0 11
	QNENAYU QNENAYU URASGRU URASGRU DNPTYIG SSGDRIYU SSGDRIYU SSGDRIYU KNANTLS KNANTLS KNANTLS KNANTLS VNSDPYI VNSDPYI DSPPQF7 DSNPQF7 DSSFQK7	HA1 HA1 SRG SRG SRG SRG SRG SRG SRG SRG SRG SRG	DGVKLES DGVKLES COUVELS COUVES
KLCRLK- KLCRLN- RLCDS BLCDS KLCDLN- KLCDLN- KLCPLN- KLCDLN- KLCPLN- KLCPS TYCSLN- TYCSLN- TYCSLN- TYCSLN- FRAUKG	KDQ0NIY KDQ0NIY REQTALY REQTALY REQTALY NEQTILY NEQTILY NEQTILY NEQTILY NEQTILY TEQTILY TEQTILY NECTILY NECTILY NECTILY NECTILY NECTILY NECTILY NECTILY	SIQ OTE OTE OTE	KLIMBERV KLIMPERV KLIMPER LUNER LUNER KLIMPER LUNER KLIMPER KL
LLEDSHNC LLESTHNC VUCESTRUE VUCESTRUE VESCTLE VESTANIC LLEVTHNC LLETHNC VUETEHNC VUETEHNC VLETHNC VLETHNC VLETHNC VLETHNC VLETHNC VLETHNC VLETHNC VLETHNC	LINER CONTRACT CONTRA	• • • • • • • • • • • • • •	YPK YSEE YPK YEEE HUY YRDEJ HUY YRDEJ HUY YRDEJ HTO YREE HTO YREE HTO YREE HTE YEEE HTE YEE HTE YEE H
TUTHSVN TUTHSVN TUTHSE EUTHAED EUTHAED EUTHAED TUTHAD TUTHSVE EUTNATE EUTNATE EUTSVE EUTSSVE PUTSSE	YFTYDSK.	RSAKLRM KGNTLKL KONTLKL KOGELVLL KOGELVLL KOSELKL RSSELKL NORELKL NORSLKL NORSLKL NORSLKL NORSLKL KSOCIKL KSOSLFL KSOSLFL KSOSLFL KSOSLFL KTP - LKL KTP - LKL	VRMGTYD IRNGTYD IRNGTYD IRNGTYD VRNGTYD VRNGTYD IRNNTYD IRNNTYD IRNGTYD IRNNTYD IRNGTYD IRNGTYD IRNGTYD IRNGTYD IRNGTYD IRNGTYD IRNGTYD IRNGTYD IRNGTYD IRNGTYD
TULEKNU TILEKNU TILEKNU TILEKNU TILEKNU TILEKNU TILEKNU TLEKNU TLERGU TLERGU TLEBGU TLEBGUU TLEBGUU TLEBGUU TLEBGUU TLEBGUU TLEBGUU TLEBGUU	EVLVLMG COMLING DKLYLWG DKLYLWG DLULWG DLULWG DLULWG DLULWG DLULWG DLULWG DULLVWG DLULVWG DULLVHUG DULLVWG DULLVLWG DULLV DULLVLWG DULLVLWG DULLVLWG DULLVLWG DULLVLWG DULLVLWG DULLVLWG DULLVLWG D	 GESCPRYUGESCERYUGESCERYUGESCERYUGESCERYUGESCERYUDESCERSCERVERUPERUUGESCERYUGESCER	DNBCMES DDBCCMES DDCCMES DDCCMES
	SYVNKKGK SYVNKKGK IMPNUGKE IMPNUGKE IMPNUGKE IMPNUGKE AYVNNGKU AYVNNEDG AYVNNEDG AYVNNEDG IVV	- IHP VTJ - IHP VTJ - IHP LTJ - INP LTJ - INP RTM - IDF	2785YHKC 2781YHKC 277
	SYFKLKN: NYELARR: KYFALNU: KYFALNU: AYSUJUL AYSUJUR AYSUJUR SYFODNI KYFUZFU KYFUZFU NYFFOLNO'	TTTLFFU TTTLFFU TTTLFFU TTTLFFU TTTLFFU TTTTLFFU TTTTLFFU TTTTLFFU TTTTLFFU TTTTTFFU TTTTTTTTTT	AKEICNG AEKELGNG AEKELGNG AEKELGNG AEKECNG AKELGNG AKEECNG AEKECNG AEKECNG AEKECNG AEKECNG AEKECNG AEKECNG AEKECNG AEKECNG AEKECNG AKEC
CCONLYAT		775LGAIN 775LGAIN 777LGAIN 777NGSIF 777NGSIA 777NGGIN 777AGGIN 777AGGIN 777AGGIN 777AGGIN 775GGIN 775SGGIN 775SGGIN 775SGGIN 775SGGIN 775SGGIN 775SGGIN 755SG	YKSQLKNN FRQLREN FRQLN
Hondrad Hon	NLLWLTE NULWLTE RLAWLARS RLAWLARS RLAWLARS RLAWLARS SERVULIS SERVULIS SERVULIS SERVULIS SERVULIS SERVULIS SERVULIS SERVULIS SERVULIS SERVULIS SERVULIS	- ECNTRC - ECNTRC - NCETRC - NCETRC - SCVSEC - SCVSEC - NCNTSC - NCNTSC - SCVSEC - SCSTSC - S	VKNLYDEY WKNLYDEY MNKLFEEK MNKLFEEK MNKLFEEK WKNLYDE VKNLYDE WNNLYKEY WKNLHDE VKNLHEY VKNLHEY WNLNLYER MNKLFEEK MLNLYER WIKLFEEK MNK MNKLFEEK MNK MNK MNK MNK MNK MNK MNK MNK MNK MN
1 IICYHANN IICYHANN IICYHANN IICYHANN IICYHANN IICYHANN IICYHANN IICYHANN IICYHAN IICYYINN IICYCHHVE IICHHVE IICHHVE IICHHVE IICHHVE IICHHVE IICHHVE IICHHVE IICHHVE IICHHVE IICHHVE IICHHVE IICHHVE IICHHVE IICH	KS - SFYR 367 - SFF 577 - SFF 577 - SFF 562 - SFF 562 - SFF 562 - SFYR 562 - SFYR 562 - SFYR 562 - SFYR 571 -	15-NASMH 17-NASMH 17-DATG 18-DAPIG 18-D	TLDFHDSN TTD7HDSN TTD7HDSN TTD7TD5 TTD7TD5 TTD7TD5 TTDLAD5E TTDLAD5E TTDLAD5E TTD2HDSN TTD8HDSN TTD8HDSN TTD7TD5 TTD7TD5 TTD7TD5 TTD7TD5 TTTD2D5 TTTC2
de HA b7117 b7117 b117 b117 b117 b017 b017 b017 b017 b017 b017 b017 b017 b017 b017 b017 b017 	AACSHA-(AACSHA-(YACKADY- YACKKA-) YACKKA-1 SACKKA-1 SACKKA- KACNKA- KACNATS KACNTSI SACKUTSI	* * * * * * * * * * * * * * * * * * *	VLLENER VLLENER VALENOR VALENOR VLANENER VLLANER VLLENOR VLLEN
pepti Advicesory Advicesory Sisservesory Advicesory Sisservesory Advicesory Advicesory Advicesory Advicesory Advicesory Advicesory Advicesory Advicesory	7	• 	WYNAELL WYNAELL WYNAELL WYNAELL SYNAELL WYNAELL MYNAELL MYNAELL MYNAELL WYNAELL WYNAELL WYNAELL WYNAELL WYNAELL WYNAELL WYNAELL WYNAELL WYNAELL WYNAELL WYNAELL WYNAELL WYNAELL WYNAELL WYNAELL WYNAELL WYNAELL WYNAELL WYNAELL
Signal LCALAAAN ILLEALANY ILLEATANY ILLEATANY ILLEATANY ILLAARN AALLANY AALLANY AALLANY AALLANY AALLANY ALLEAVON ILLESVON ILLESVON ILLESVON ILLESVON ILLESVON ILLESVON	WENHNTTI WTOLITTY WTOLITTY WTOLITTY WTOLITTY WINTTS WINTS WINTTS	YAFALSR(YYERISK) GYFXLINY GYFXLINY GYFXLINY GYFXLINY YAFLVK) YAFLUFU YGHVLFO YGHVLFO YGHUTO YGHLIFO GYFILEE G	HIGHIGAN HIGHIG
L MKANLLVL MKANLLVL MKANLLVL MKANLLVL MKENTIAL MKENTIAL MKENTLF MKENTLF MKENTLF MKALTLV MKATLLTV MKATLLTV MKATLLLV	FELTEPRESS FELTEPRESS FELTEPRESS FELTEPRESS FELTEPRESS CALLER FELTEPS FELTEPRESS FELTEPRESS FELTEPS FE	161 ANGNLIZAF STGNLIZAF STGNLIZAF STGNLIZAF STGNLIZAF STGNLIZAF STGNLIZAF NGGATAF NGGATAF NGGATAF NGGALIAF STGGLIZAF	541 REENLARKK REENLARKK REDELBRUKKK REODLERKY REODLERKY REDDLERKY
C 8 12 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4	нн нн нн нн нн нн нн нн с нн нн с с с с	C B 111211 C B 1142 C	НН Н Н Н Н Н Н Н Н Н Н Н Н Н Н С С В В Н Н 12 С Н Н 8 С Н Н 12 С Н Н 12 С С Н 12 С Н 12 С Н 12 С Н 12 С Н 12 С Н 12 С Н 12 С Н 12 С Н 12 С Н 12 С Н 12 С Н 12 С Н 12 С Н 12 С С Н 12 С Н 12 С Н 12 С Н 12 С Н 12 С Н 12 С Н 12 С Н 12 С Н 12 С Н 12 С Н 12 С Н 12 С Н 12 С Н 12 С Н 12 С Н 12 С Н 12 С Н 12 С С Н 12 С С Н 12 С С Н 12 С С Н 12 С С Н 12 С С Н 12 С С Н 12 С С Н 12 С С С С С С С С С С С С С С С С С С

Table 1.	
Amino Acid Sequences of the HA or HE Genes used for a Phylogeneic Analysis of A	
Viruses (15 subtypes) and B and C Viruses ^a	

Type and Subtype	Virus ^b	Accession Number ^c
Influenza A virus H1	A/PR/8/34 (H1N1)	V01088
Influenza A virus H2	A/herring gull/DE/677/88 (H2N8)	L11132
Influenza A virus H3	A/Beijing/353/89 (H3N2)	U97740
Influenza A virus H4	A/budgerigar/Hokkaido/1/77 (H4N6)	M25285
Influenza A virus H5	A/goose/Guangdong/1/96 (H5N1)	AF144305
Influenza A virus H6	A/shearwater/Australia/1/72 (H6N5)	D90303
Influenza A virus H7	A/turkey/Oregon/71 (H7N3)	M31689
Influenza A virus H8	A/turkey/Ontario/6118/68 (H8N4)	D90304
Influenza A virus H9	A/turkey/Wisconsin/66 (H9N2)	D90305
Influenza A virus H10	A/chicken/Germany/N/49 (H10N7)	M21647
Influenza A virus H11	A/duck/England/56 (H11N6)	D90306
Influenza A virus H12	A/duck/Alberta/60/76 (H12N5)	D90307
Influenza A virus H13	A/gull/Maryland/704/77 (H13N6)	D90308
Influenza A virus H14	A/mallard duck/Gurjev/263/82 (H14N5)	M35997
Influenza A virus H15	A/duck/Australia/341/83 (H15N8)	L43916
Influenza B virus	B/Lee/40	K00423
Influenza C virus	C/Johannesburg/66	M17868

^a In the phylogenetic analysis, only the HA2 (HE2) region of the HA (HE) gene was used.

^b The virus name consists of type/host species/place of isolation/strain number if available/year of isolation (the HA subtype and the neuraminidase subtype). By convention, the host species is omitted if the virus is isolated from humans (WHO Memorandum 1980).

^c The accession number in the international DNA databank (DDBJ/EMBL/GenBank).

DNA databank (DDBJ release 43). After excluding sequences from laboratory-adapted viruses and identical sequences within species, we obtained 57, 34, 58, 10, 29, 2, 41, 1, 4, 2, 1, 1, 3, 1, and 2 amino acid sequences for the H1–H15 subtypes of A virus HA2s, respectively. We also obtained 15 sequences for B virus HA2s and 35 sequences for C virus HE2s. A total of 296 amino acid sequences were aligned by the computer program CLUSTAL W (Thompson, Higgins, and Gibson 1994). After removing all alignment gaps, 207 amino acid sites were used for estimating p, Poisson correction (PC), and gamma distances (Nei and Kumar 2000). The gamma shape parameter (a) was estimated to be 1.83 by Gu and Zhang's (1997) method. The phylogenetic tree was constructed by the neighbor-joining (NJ) method (Saitou and Nei 1987), and the reliability of each interior branch was tested by the bootstrap method with 1,000 resamplings (Felsenstein 1985; Kumar et al. 2001). The NJ trees were also constructed for 17 amino acid sequences which were randomly chosen from each subtype of A virus HA2s and from B virus HA2s and C virus HE2s (table 1).

Estimation of Divergence Times

For estimating the divergence times between subtypes of influenza A virus genes, we used only A virus sequences because B and C virus sequences were not necessary. We also used amino acid sequences for the entire region of HA because the alignment for A virus HAs appeared to be reliable (fig. 1) (Rohm et al. 1996), and longer sequences were expected to give more reliable estimates.

We obtained 50, 25, 24, 10, 21, 2, 25, 1, 4, 2, 1, 1, 3, 1, and 2 amino acid sequences for the H1–H15 subtypes of A virus HAs from the databank, respective-

ly, and made a multiple alignment for a total of 172 sequences by CLUSTAL W. After removing all alignment gaps, 540 amino acid sites were used for estimating gamma distances with a = 1.20, which was obtained by Gu and Zhang's method. An NJ tree was constructed, and the branch lengths were recalculated by the ordinary least squares method (Rzhetsky and Nei 1993) to estimate the rate of amino acid substitution accurately (see subsequently).

When the years of isolation are available for viral sequences in a phylogenetic tree, the rate of amino acid substitution may be estimated by the regression coefficient of the numbers of amino acid substitutions from a common root on the years of isolation (Nei 1983; Suzuki, Wyndham, and Gojobori 2001). Using the phylogenetic tree for 172 sequences of influenza A virus HAs, we estimated the rate of amino acid substitution for duck A virus HAs because duck provided the largest number (28) of sequences among aquatic birds. For estimating the divergence times between subtypes of A virus HA genes, we constructed a linearized tree (Takezaki, Rzhetsky, and Nei 1995) for 28 amino acid sequences of duck A virus HAs using the gamma distance with a = 1.20. The standard errors (SEs) and 99% confidence intervals (CIs) of the rates and the divergence times were estimated by the bootstrap method, under the assumption that the topologies of the phylogenetic trees for 172 sequences of influenza A virus HAs and 28 sequences of duck A virus HAs were correct (Nei and Kumar 2000).

Results

Phylogenetic Relationships of Influenza A, B, and C Virus HA Genes

The NJ trees constructed by using p, PC, and gamma distances for 17 randomly chosen amino acid se-



FIG. 2.—Phylogenetic trees for 17 randomly chosen amino acid sequences (see table 1) of the HA2 (HE2) region of influenza A, B, and C virus HAs (HEs). Trees (a), (b), and (c) were obtained by using p, PC, and gamma distances, respectively. The bootstrap value is indicated for each interior branch. The scale bars indicate the numbers of amino acid substitutions per site.

quences of the HA2 (HE2) protein are shown in panels (a), (b), and (c) of figure 2, respectively. All trees show the same topology and indicate that all influenza A virus HA genes diverged after they separated from B virus HA genes. The monophyly of A virus HA genes is supported by a bootstrap value of 100%, 99%, and 95% in trees (a), (b), and (c), respectively. This relationship was also supported by the NJ trees for 296 amino acid sequences of influenza A, B, and C virus HA2s (HE2s) with high bootstrap values (100%, 99%, and 86% for p, PC, and gamma distances, respectively) (data not shown).

Divergence Times Between Subtypes of A Virus HA Genes

For estimating the divergence times between subtypes of A virus HA genes, we first estimated the rate of amino acid substitution for duck A virus HAs because duck is one of the natural reservoirs of these viruses and provided the largest number of sequences among them. In the phylogenetic tree for 172 amino acid sequences of the entire region of A virus HAs, only the H1 and H2 subtypes included sufficient numbers of sequences for estimating the rate for duck A virus HAs and are



FIG. 3.—(a) H1 and (b) H2 subtype subtrees in the phylogenetic tree for 172 amino acid sequences of the entire region of influenza A virus HAs. For the nomenclature of influenza viruses, see the footnotes of table 1. The nodes at which the year of divergence was estimated (table 3) are denoted by black circles and labeled with letters. The human and swine sequences used for estimating the rates and the years of divergences are underlined. The duck sequences used for estimating the rates are boldfaced. The scale bars indicate the numbers of amino acid substitutions per site.

shown in panels (a) and (b) of figure 3, respectively. In this figure, avian sequences had generally shorter branch lengths than human and swine sequences in both subtypes, indicating that the rate for the former was slower than that for the latter. To estimate the rate of amino acid substitution, we conducted a regression analysis using duck sequences but failed to obtain the rate because it became negative in both H1 and H2 subtypes (data not shown). This happened probably because the evolutionary rate for duck sequences was too slow to give reliable estimates (Bean et al. 1992; Schafer et al. 1993; Suarez 2000).

To obtain a reliable rate for duck sequences, it was necessary to analyze duck sequences which were more distantly related from one another than those analyzed earlier in the article. For this purpose, we estimated the



FIG. 4.—Regression analyses for estimating the rate of amino acid substitution for the entire region of (*a*) human, (*b*) classical swine, and (*c*) avian-like swine influenza A virus HAs in the H1 subtype and for (*d*) human sequences in the H2 subtype and for the (*e*) H1 and (*f*) H2 subtype sequences of duck HAs. In each panel, the abscissa indicates the year of isolation or divergence, and the ordinate indicates the number of amino acid substitutions from node (a) M, (b) M, (c) N, (d) O, or (e) M, or from (f) the earliest node of panel (b) in figure 3. An open circle and an open square in panel (e) indicate nodes M and N in figure 3*a*, respectively, and an open triangle in panel (f) indicates node O in figure 3*b*.

years of divergences between duck sequences and human and swine sequences in figure 3, using the rates for the latter sequences and added these nodes to the regression analysis of duck sequences. The rates for human and swine sequences were easily estimated by the regression analysis using these sequences only because

Table 2. Rates of Amino Acid Substitution (× 10^{-3}) for Influenza A Virus HAs^a

Subtype	Host Species	Rate ± SE (Per Site Per Year)	99% CI (Per Site Per Year)
H1	Human	$\begin{array}{c} 1.20 \pm 0.32 \\ 0.56 \pm 0.20 \\ 1.87 \pm 0.79 \\ 0.39 \pm 0.12 \\ 2.03 \pm 0.77 \\ 0.25 \pm 0.10 \end{array}$	0.49-2.11
H1	Classical swine		0.16-1.21
H1	Avian-like swine		0.30-4.21
H1	Duck		0.14-0.72
H2	Human		0.15-4.01
H2	Duck	0.25 ± 0.19	-0.23-0.85
H1-H2 ^b		0.32 ± 0.11	0.07-0.66

^a The entire region of the HA gene was used for analysis.

^b The average rate for the H1 and H2 subtypes.

the rates were relatively high. In the H1 subtype, we first estimated the year of divergence at node M using the rate for human sequences (fig. 3a). We used only human sequences which were isolated before 1977 because human A viruses circulating after 1977 are known to have originated from a laboratory-adapted virus (Kendal et al. 1978; Nakajima, Desselberger, and Palese 1978; Scholtissek, von Hoyningen, and Rott 1978; Palese and Young 1982; Hayashida et al. 1985). The rate of amino acid substitution for human sequences was estimated to be 1.20×10^{-3} per site per year (fig. 4a and table 2), and the year of divergence at node M was A.D. 1862. We also estimated the year of divergence at the same node using classical swine sequences. This gave a rate of 0.56×10^{-3} per site per year (fig. 4b) and a date of 1836 at node M. Taking the average for human and classical swine sequences, we obtained 1849 as the final estimate of the year of divergence at node M (table 3). In addition, we estimated the year of divergence at node N using avian-like swine sequences. These sequences apparently evolved at a rate of 1.87×10^{-3} per site per year (fig. 4c), and the year of divergence at node N was

Table 3.Years of Divergences Between Influenza A Virus HAGenes in the H1 and H2 Subtypes

Y Node ^a	tear of Divergence \pm SE (A.D.)	99% CI (a.d.)
M	$ 1849 \pm 32 1965 \pm 53 1946 \pm 80 $	1660–1886 1843–1977 1781–1953

^a The nodes correspond to those indicated in figure 3.

1965. By adding nodes M and N to the regression analysis of duck sequences (fig. 4e), the rate for these sequences was estimated to be 3.89 \times 10⁻⁴ per site per year. We also estimated the rate for duck sequences in the H2 subtype. In this subtype, we estimated the year of divergence at node O using human sequences (fig. 3b). The rate of amino acid substitution for these sequences was estimated to be 2.03×10^{-3} per site per year (fig. 4d and table 2) and the year of divergence at node O was 1946 (table 3). We added node O to the regression analysis of duck sequences (fig. 4f) and obtained a rate of 2.48×10^{-4} per site per year for these sequences. Taking the average for the H1 and H2 subtypes, we obtained 3.19×10^{-4} per site per year as the final estimate of the rate of amino acid substitution for duck influenza A virus HAs.

A linearized tree for 28 amino acid sequences of duck A virus HAs is shown in figure 5. The topology of subtypes was the same as that shown in figure 1, except for the branching pattern of the H11 subtype, which was supposed to make a cluster with the H12 subtype but made a cluster with the H1, H2, and H5 subtypes. The estimates of the divergence times between different subtypes of influenza A virus HA genes are listed in table 4. Although the SE and 99% CI are large, all subtypes apparently diverged from several thousand to several hundred years ago. In particular, the earliest divergence (node X) is likely to have occurred about 2,000 years ago. We further estimated the divergence times between influenza A, B, and C virus HA (HE) genes by linearizing the phylogenetic tree in figure 2c. Assuming that the earliest divergence between subtypes of A virus HA genes occurred 1,971 years ago (table 4), A and B virus HA genes apparently diverged 3,832 years ago, and the separation of A and B virus HA genes from C virus HE genes occurred 7,919 years ago.

Discussion

The divergence between influenza A and B virus HA genes apparently occurred earlier than the divergences between different subtypes of A virus HA genes. This is different from the conclusion of Webster et al. (1992) that the divergence between A and B virus HA genes occurred later than the divergences between some subtypes of A virus HA genes. This difference is apparently caused by the fact that Webster et al. did not use an outgroup for their phylogenetic analysis. In the present study, the phylogenetic relationship was supported with a high bootstrap value, indicating that the relationship obtained is highly reliable.

The rate of amino acid substitution for duck A virus HAs $(3.19 \times 10^{-4} \text{ per site per year})$ was slower than that for human and swine A virus HAs ($[0.56-2.03] \times 10^{-3}$ per site per year) but similar to that for B virus HAs $(5.3 \times 10^{-4} \text{ per site per year [Air et al. 1990]})$ and C virus HEs $(2.3 \times 10^{-4} \text{ per site per year [Muraki et$ $al. 1996]})$. These results suggest that the rate for HAs (HEs) is more or less constant in the natural reservoir but is accelerated in the newly infected host species. This is probably caused by variation in the strengths of immune responses and functional constraints on HAs (HEs) among different host species (Yamashita et al. 1988; Bean et al. 1992; Schafer et al. 1993; Scholtissek, Ludwig, and Fitch 1993; Makarova et al. 1999; Suzuki and Gojobori 1999).

The earliest divergence time between subtypes of influenza A virus HA genes was estimated to be about 2,000 years ago. Also, the divergence time between A and B virus HA genes was estimated to be about 4,000 years ago, whereas A and B virus HA genes and C virus HE genes diverged about 8,000 years ago. These estimates are substantially higher than those (200-300 years) by Saitou and Nei (1986), who used human HA sequences. Because the evolutionary rate for human A virus HAs is known to be higher than that for aquatic birds, their estimates are considered to be underestimates. In fact, influenza pandemics in humans have been recorded as early as 412 B.C. (Kaplan and Webster 1977), suggesting that influenza A viruses existed more than 2,400 years ago. This observation is consistent with the estimates obtained in the present study.

We estimated the rates and the divergence times under the assumption that the molecular clock has held throughout the evolutionary history of HA (HE) genes. To examine whether this was really the case, we tested the linear relationship between the year of isolation and the number of amino acid substitutions in figure 4 and found that the linearity was not supported at the 1% significance level in both panels (a) and (f). However, the rate of amino acid substitution for human A virus HAs obtained from panel (a) $(1.20 \times 10^{-3} \text{ per site per})$ year) was similar to that from previous studies (1.0 \times 10^{-3} per site per year [Saitou and Nei 1986]), and the rate for duck A virus HAs obtained from panel (f) (3.89 \times 10⁻⁴ per site per year) was similar to that obtained from panel (e) $(2.48 \times 10^{-4} \text{ per site per year})$. These observations suggest that the rates obtained from panels (a) and (f) are approximately correct. Also, the molecular clock was not rejected at the 1% significance level for the phylogenetic tree in figure 5 by the likelihoodratio test (Rambaut 2000; Yang 2000) but was rejected for the tree in figure 2c. The latter observation may reflect the fact that the biochemical functions are different between HAs and HEs and the natural reservoirs are not the same for influenza A, B, and C viruses. Therefore, some caution is necessary in estimating the divergence times between influenza A, B, and C virus HA (HE) genes. However, the rate of amino acid substitution for duck influenza A virus HAs was similar to that for B virus HAs and C virus HEs, as indicated previously. Also, in reality, no strict molecular clock is likely to



FIG. 5.—Linearized tree for 28 amino acid sequences of the entire region of duck influenza A virus HAs. For the nomenclature of influenza viruses, see the footnotes of table 1. The nodes at which the divergence times were estimated (table 4) are denoted by black circles and labeled with letters. The scale bar indicates the number of years before present (above the line) and the number of amino acid substitutions per site (below the line).

hold for any protein but it is known that rough divergence times can be obtained even if the molecular clock is violated to some extent (Nei and Kumar 2000, pp. 187–206; Nei, Xu, and Glazko 2001). Therefore, these estimates also appear to be appropriate as rough estimates.

In conclusion, influenza virus HA (HE) genes apparently evolved at a rate of amino acid substitution of 10^{-4} per site per year in the natural reservoir. These

Table 4. Divergences Times Between Subtypes of Influenza A Virus HA Genes

Nodea	Divergence time ± SE (years ago)	99% CI (years ago)
P	443 ± 43	342–579
Q	655 ± 57	519-819
R	$1,175 \pm 90$	966–1,439
<u>S</u>	$1,533 \pm 119$	1,252–1,870
Τ	379 ± 41	279-491
U	433 ± 45	331–562
V	722 ± 60	576-882
W	$1,433 \pm 105$	1,204–1,731
Χ	$1,971 \pm 132$	1,680–2,336

^a The nodes correspond to those indicated in figure 5.

genes apparently diverged into influenza A, B, and C virus HA (HE) genes several thousand of years ago and subsequently into subtypes in influenza A viruses from several thousand to several hundred years ago.

Acknowledgments

The authors thank two anonymous reviewers for their valuable comments. This study was supported by grants from the National Institutes of Health to M.N. (GM20293). Y.S. is supported by the JSPS Research Fellowships for Young Scientists.

LITERATURE CITED

- AIR, G. M., A. J. GIBBS, W. G. LAVER, and R. G. WEBSTER. 1990. Evolutionary changes in influenza B are not primarily governed by antibody selection. Proc. Natl. Acad. Sci. USA 87:3884–3888.
- BEAN, W. J., M. SCHELL, J. KATZ, Y. KAWAOKA, C. NAEVE, O. GORMAN, and R. G. WEBSTER. 1992. Evolution of the H3 influenza virus hemagglutinin from human and nonhuman hosts. J. Virol. 66:1129–1138.
- COX, N. J., F. FULLER. N. KAVERIN, H. D. KLENK, R. A. LAMB, B. W. J. MAHY, J. MCCAULEY, K. NAKAMURA, P. PALESE, and R. WEBSTER. 2000. Family *Orthomyxoviridae*. Pp. 585–

597 *in* M. H. V. VAN REGENMORTEL, C. M. FAUQUET, D. H. L. BISHOP, et al. (11 co-editors), eds. Virus Taxonomy. Academic Press, London.

- FELSENSTEIN, J. 1985. Confidence limits on phylogenies: an approach using the bootstrap. Evolution **39**:783–791.
- GAMMELIN, M., A. ALTMULLER, U. REINHARDT, J. MANDLER, V. R. HARLEY, P. J. HUDSON, W. M. FITCH, and C. SCHOL-TISSEK. 1990. Phylogenetic analysis of nucleoproteins suggests that human influenza A viruses emerged from a 19thcentury avian ancestor. Mol. Biol. Evol. 7:194–200.
- GU, X., and J. ZHANG. 1997. A simple method for estimating the parameter of substitution rate variation among sites. Mol. Biol. Evol. 14:1106–1113.
- HAYASHIDA, H., H. TOH, R. KIKUNO, and T. MIYATA. 1985. Evolution of influenza virus genes. Mol. Biol. Evol. 2:289– 303.
- HINSHAW, V. S., R. G. WEBSTER, and B. TURNER. 1980. The perpetuation of orthomyxoviruses and paramyxoviruses in Canadian waterfowl. Can. J. Microbiol. **26**:622–629.
- KAPLAN, M. M., and R. G. WEBSTER. 1977. The epidemiology of influenza. Sci. Am. 12:88–106.
- KENDAL, A. P., G. R. NOBLE, J. J. SKEHEL, and W. R. DOWDLE. 1978. Antigenic similarity of influenza A (H1N1) viruses from epidemics in 1977–1978 to "Scandinavian" strains isolated in epidemics of 1950–1951. Virology 89:632–636.
- KROSSOY, B., I. HORDVIK, F. NILSEN, A. NYLUND, and C. EN-DRESEN. 1999. The putative polymerase sequence of infectious salmon anemia virus suggests a new genus within the *Orthomyxoviridae*. J. Virol. **73**:2136–2142.
- KUMAR, S., K. TAMURA, I. B. JAKOBSEN, and M. NEI. 2001. MEGA2: molecular evolutionary genetics analysis software. Bioinformatics 17:1244–1245.
- MAKAROVA, N. V., N. V. KAVERIN, S. KRAUSS, D. SENNE, and R. G. WEBSTER. 1999. Transmission of Eurasian avian H2 influenza virus to shorebirds in North America. J. Gen. Virol. 80:3167–3171.
- MURAKI, Y., S. HONGO, K. SUGAWARA, F. KITAME, and K. NAKAMURA. 1996. Evolution of the haemagglutinin-esterase gene of influenza C virus. J. Gen. Virol. 77:673–679.
- NAKADA, S., R. S. CREAGER, M. KRYSTAL, R. P. AARONSON, and P. PALESE. 1984. Influenza C virus hemagglutinin: comparison with influenza A and B virus hemagglutinins. J. Virol. **50**:118–124.
- NAKAJIMA, K., U. DESSELBERGER, and P. PALESE. 1978. Recent human influenza A (H1N1) viruses are closely related genetically to strains isolated in 1950. Nature 274:334–339.
- NEI, M. 1983. Genetic polymorphism and the role of mutation in evolution. Pp. 165–190 *in* M. NEI and R. K. KOEHN, eds. Evolution of genes and proteins. Sinauer, Sunderland, Mass.
- NEI, M., and S. KUMAR. 2000. Molecular evolution and phylogenetics. Oxford University Press, Oxford, New York.
- NEI, M., P. XU, and G. GLAZKO. 2001. Estimation of divergence times from multiprotein sequences for a few mammalian species and several distantly related organisms. Proc. Natl. Acad. Sci. USA 98:2497–2502.
- PALESE, P., and J. F. YOUNG. 1982. Variation of influenza A, B, and C viruses. Science **215**:1468–1474.
- RAMBAUT, A. 2000. Estimating the rate of molecular evolution: incorporating noncontemporaneous sequences into maximum likelihood phylogenies. Bioinformatics **16**:395–399.
- REID, A. H., T. G. FANNING, J. V. HULTIN, and J. K. TAUBEN-BERGER. 1999. Origin and evolution of the 1918 "Spanish" influenza hemagglutinin gene. Proc. Natl. Acad. Sci. USA 96:1651–1656.

- ROHM, C., N. ZHOU, J. SUSS, J. MACKENZIE, and R. G. WEB-STER. 1996. Characterization of a novel influenza hemagglutinin, H15: criteria for determination of influenza A subtypes. Virology 217:508–516.
- RZHETSKY, A., and M. NEI. 1993. Theoretical foundation of the minimum-evolution method of phylogenetic inference. Mol. Biol. Evol. 10:1073–1095.
- SAITOU, N., and M. NEI. 1986. Polymorphism and evolution of influenza A virus genes. Mol. Biol. Evol. **3**:57–74.
- ——. 1987. The neighbor-joining method: a new method for reconstructing phylogenetic trees. Mol. Biol. Evol. 4:406– 425.
- SCHAFER, J. R., Y. KAWAOKA, W. J. BEAN, J. SUSS, D. SENNE, and R. G. WEBSTER. 1993. Origin of the pandemic 1957 H2 influenza A virus and the persistence of its possible progenitors in the avian reservoir. Virology **194**:781–788.
- SCHOLTISSEK, C., S. LUDWIG, and W. M. FITCH. 1993. Analysis of influenza A virus nucleoproteins for the assessment of molecular genetic mechanisms leading to new phylogenetic virus lineages. Arch. Virol. 131:237–250.
- SCHOLTISSEK, C. V., V. VON HOYNINGEN, and R. ROTT. 1978. Genetic relatedness between the new 1977 epidemic strains (H1N1) of influenza and human influenza strains isolated between 1947 and 1957 (H1N1). Virology 89:613–617.
- SLEMONS, R. D., D. C. JOHNSON, J. S. OSBORN, and F. HAYES. 1974. Type-A influenza viruses isolated from wild free-flying ducks in California. Avian Dis. 18:119–124.
- SMITH, W., C. H. ANDREWES, and P. P. LAIDLAW. 1933. A virus obtained from influenza patients. Lancet **225**:66–68.
- SUAREZ, D. L. 2000. Evolution of avian influenza viruses. Vet. Microbiol. 74:15–27.
- SUZUKI, Y., and T. GOJOBORI. 1999. A method for detecting positive selection at single amino acid sites. Mol. Biol. Evol. **16**:1315–1328.
- SUZUKI, Y., A. WYNDHAM, and T. GOJOBORI. 2001. Virus evolution. Pp. 377–413 in D. J. BALDING, M. BISHOP, and C. CANNINGS, eds. Handbook of statistical genetics. Wiley, Chichester.
- TAKEZAKI, N., A. RZHETSKY, and M. NEI. 1995. Phylogenetic test of the molecular clock and linearized trees. Mol. Biol. Evol. **12**:823–833.
- 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.
- WEBSTER, R. G., W. J. BEAN, O. T. GORMAN, T. M. CHAMBERS, and Y. KAWAOKA. 1992. Evolution and ecology of influenza A viruses. Microbiol. Rev. 56:152–179.
- WEBSTER, R. G., M. YAKHNO, V. S. HINSHAW, W. J. BEAN, and K. G. MURTI. 1978. Intestinal influenza: replication and characterization of influenza viruses in ducks. Virology 84: 268–278.
- WHO MEMORANDUM. 1980. A revision of the system of nomenclature for influenza viruses. Bull. WHO 58:585–591.
- YAMASHITA, M., M. KRYSTAL, W. M. FITCH, and P. PALESE. 1988. Influenza B virus evolution: co-circulating lineages and comparison of evolutionary pattern with those of influenza A and C viruses. Virology 163:112–122.
- YANG, Z. 2000. Phylogenetic analysis by maximum likelihood (PAML). Version 3.0. University College London, London, U.K.

NARUYA SAITOU, reviewing editor

Accepted December 4, 2001