

乙型肝炎病毒蛋白表型定义初探

董菁, 任建林, 卢雅丕

董菁, 任建林, 卢雅丕, 厦门大学医学院第一临床学院消化内科 福建省厦门市 361004

董菁, 男, 1969年2月生人, 河北省徐水县人, 汉族. 2001年北京大学医学部毕业, 医学博士, 主治医师. 主要从事乙型肝炎病毒分子生物学与病毒性肝炎的治疗的研究.

福建省卫生厅青年科研课题资助: No. 2004-1-26

项目负责人: 董菁, 361004, 福建省厦门市, 厦门大学医学院第一临床学院消化内科. dj@xmzsh.com

电话: 0592-2292017 传真: 0592-2292017

收稿日期: 2004-05-07 接受日期: 2004-06-17

Definition of prototype of hepatitis B virus: A preliminary study

Jing Dong, Jian-Lin Ren, Ya-Pi Lu

Jing Dong, Jian-Lin Ren, Ya-Pi Lu, Department of Gastroenterology, the First Clinical College of Xiamen University, Amoy 361004, China. Supported by the Science and Technology Foundation of Fujian Province for Young Scholars, No. 2004-1-26

Correspondence to: Jing Dong, Department of Gastroenterology, the First Clinical College of Xiamen University, Amoy 361004, China. dj@xmzsh.com

Received: 2004-05-07 Accepted: 2004-06-17

Abstract

AIM: To create a new typing method showing the difference among HBV strains after reviewing the HBV genome sequences labeled with different genotypes in the GenBank.

METHODS: HBV genome sequences were collected from the GenBank and then classified into 8 groups based on their genotypes labeled by authors. The Vector NTI suite 8.0 software was used to compare the identity and difference among the strains of HBV genomes. Possible regions encoding pre-pre-S, pre-X and pre-C peptides were also analyzed with this software.

RESULTS: One hundred and nineteen full-length HBV genomes from GenBank were collected, and then sorted into 8 groups according to their genotypes. The total positive rate and total identical rate of 119 sequences were 95.7% and 47.7%, respectively. The total positive rates of whole C protein, whole S protein, whole X protein and polymerase amino acids sequences were 98.6%, 87.3%, 57.2% and 95.2%, respectively; and total identical rates were 37.4%, 24.1%, 27.7% and 43.5%. In the study group, 33.61% strains encoded pre-pre-S peptide, 14.3% strains encoded pre-X peptide, 26.05% strains had no function of encoding pre-C peptide, whereas, 94.1% of pre-X coding strains also encoded pre-pre-S peptide. The identical rates of region 1-700 nt and 1 103-1 653 nt of HBV genome were 30.6% and 20.8%, respectively, and therefore they were considered as hypervariable region; the identical rate of region 1 654-1 950 nt of HBV genome was 74.2% and defined as hyperconversible region. Hypervariable and hyperconversible regions could be found in all of the four viral proteins. Based on mutations of leading peptides of

the three HBV viral proteins, a novel typing method named prototype was therefore generated. In this new category, 7 prototypes were listed, and there were 39.5% strains belonging to the major one, type IV, type V and type both covering 19.3%. All 7 prototypes were found in Asia, with the percentages of I, IV, V and VII types above 20%. There was no strain isolated from Europe belonging to prototypes I, II, or III, and the percentages of IV was 58.3%, V 13.9% and VII 25.0%, respectively.

CONCLUSION: Hypervariable and hyperconversible regions are noticed while analyzing HBV genome sequences. Furthermore, prototype, a novel term is raised to elucidate encoding of the 3 leading peptides and structural variation of viral proteins due to gene mutation.

Dong J, Ren JL, Lu YP. Definition of prototype of hepatitis B virus: A preliminary study. *Shijie Huaren Xiaohua Zazhi* 2004;12(9):2074-2085

摘要

目的: 探讨乙型肝炎病毒(HBV)基因型的分型方式, 并根据病毒蛋白结构提出新的病毒蛋白分型方式.

方法: 自GenBank中按基因型搜索符合要求的HBV基因组序列, 并应用Vector NTI suite 8.0版软件进行基因组核苷酸及各基因编码蛋白质序列比较, 并利用软件分析前-S基因、前-X基因和前-C基因的存在状态.

结果: 在GenBank中根据HBV基因型分型搜索出119个病毒株全基因组, 比较后发现选择病毒株基因组核苷酸序列总阳性率和总一致率分别为95.7%和47.7%; 选择病毒株编码的全C蛋白、全S蛋白、全X蛋白和多聚酶的总阳性率分别为98.6%、87.3%、57.2%和95.2%, 总一致率分别为37.4%、24.1%、27.7%和43.5%. 在病毒群中, 33.61%的病毒株编码前-S多肽, 14.3%的病毒株编码前-X多肽, 26.1%的病毒株不编码前-C多肽, 94.1%编码前-X多肽的病毒株同时编码前-S多肽. 基因组1-700 nt一致率30.6%, 1 103-1 653 nt一致率20.8%, 为高变区; 基因组1 654-1 950 nt的一致率为74.2%, 为高保守区. 4种病毒蛋白各有其相应的高变区和高保守区. 根据病毒蛋白前导性序列的变异情况提出新的分型方法, 命名为蛋白表型. 蛋白表型分7型, IV型为主要流行表型, 占39.5%, V型和VII型各占19.3%. 亚洲HBV蛋白分型分布分散, I、IV、V和VII型所占比例均大于20%; 欧洲IV型占58.3%, VII型占25.0%, V型占13.9%.

结论: 在综合分析HBV基因组的基础上, 初步划分出HBV

基因组和病毒蛋白内部存在的高变区和高保守区.提出蛋白表型的新概念,并综合展示基因核苷酸突变所导致的病毒蛋白的结构差异.

董菁,任建林,卢雅杰.乙型肝炎病毒蛋白表型定义初探.世界华人消化杂志 2004;12(9):2074-2085

<http://www.wjgnet.com/1009-3079/12/2074.asp>

0 引言

1968年发现了乙型肝炎病毒(HBV)的抗原,1972年Le Bouvier *et al*^[1]提出HBV表面抗原(HBsAg)根据血清反应的不同而分为不同亚型,即提出血清型分型概念. Galibert *et al*^[2]于1979年第1次解读了HBV *ayw*血清型基因组的核苷酸序列,长度为3 182 nt. 1988年Okamoto *et al*^[3]首次提出HBV基因型的概念,即根据各基因组之间差异大于8%而人为的将病毒群划分为不同的基因型,1990年Norder *et al*^[4]基于多聚酶链反应(PCR)方法建立了简便的基因型别分析方法,之后学者在研究中不断提出存在新的基因型.目前将HBV分为8种基因型,分别为A、B、C、D、E、F、G、H型^[5],基因型分布具有一定的地理特征,A型主要分布在北欧、西欧和北美,B和C型流行于东亚和远东,D型分布广泛,在地中海、印度、近东和中东地区多见,E型流行于西撒哈拉地区,F型主要在美洲大陆流行,G型主要在美国^[6],H型在中美洲流行^[7].国内主要的基因型为B、C两型^[8-9],台湾学者报告提示除E型外,其他基因型均可在华人HBV感染者中被检出,B和C型占患者人群的85%^[10].我们早期的研究认为不同的HBV病毒株的4个开放读码框架(ORF)分区中存在差异^[11],部分病毒分离株的S基因在原有的前S1区之前,存在有前前-S区编码前前-S多肽^[12];而X基因之前可能存在前-X区,编码前-X多肽^[13].在本组早期的研究中认为前前-S区和前-X区的存在可能具有基因型特异性,提出应当对中国HBV流行株的结构与功能复杂性进行重新认识^[14-15],本研究分析了目前存储在GenBank中的不同基因型HBV病毒株基因组,探讨了一种新的HBV分型方法.

1 材料和方法

1.1 材料 应用生物信息学技术进行研究,利用的材料为GenBank中存储的HBV全基因组序列.

1.2 方法

1.2.1 病毒株的选择 进入美国国立卫生院(NIH)网站,在GenBank中搜寻HBV基因组序列,之后进一步限定基因型分型分别为A、B、C、D、E、F、G、H型,将搜寻结果下载以进一步分析,其中也包括本研究组以往报告^[16-17]的全基因组序列.所筛选出的序列首先检查其病毒蛋白序列的完整性,部分病毒株虽然标明基因

型别,但其序列不能编码一种或一种以上HBV病毒蛋白,或序列中含有少量测定不准确的核苷酸位点(r或n等),这些病毒株序列被排除在本研究之外.

1.2.2 核苷酸序列分析 应用Vector 8.0版软件对下载的存储于GenBank中不同基因型的HBV基因组序列进行比较.比较前将原存储于GenBank中的HBV基因组序列进行了起始点的统一计数处理,即与Gunther *et al*^[18]和董菁 *et al*^[19]文献中HBV基因组序列的排列方式一致,各序列均以5'-TTT TTC ACC TCT GC-3'为开始,保证了各序列之间的可比性.

1.2.3 编码病毒蛋白氨基酸序列分析 将各病毒株基因组报告者所公布的P基因和前-C/C基因编码产物,即多聚酶和HBeAg氨基酸序列收集后,应用Vector 8.0版软件进行比较序列之间的一致率.利用Vector 8.0版软件具有的ORF判读功能,重新判读X、S基因,判断各病毒株是否编码前-X区^[13,20]和前前-S区^[12,21].将表达前-C多肽和核心蛋白的病毒蛋白命名为全C蛋白;如病毒株存在前-X基因和前前-S基因序列,将核酸序列进行翻译后获得的前-X多肽和前前-S多肽与原X蛋白和原S蛋白命名为全X蛋白和全S蛋白,将来自不同基因型的HBV全X蛋白(原X蛋白)和全S蛋白(原S蛋白)进行比较. Vector 8.0版软件将所有选择出的HBV序列进行比较后,提供以下重要数据:一致性序列,是软件自动比较所有序列,参考每个对应的核苷酸/蛋白质位点上不同克隆的编码/表达方式,由软件形成最具代表性的一致性序列;如果某位点出现3种以上的编码/表达方式,一致性序列中提示为空缺,表明了该位点的核苷酸/氨基酸多样性.阳性率是选定区域一致性序列的核苷酸/氨基酸序列数目与区域核苷酸/蛋白质序列总长度之比,提示差异位点的比例,用于展示区域内部的插入突变/缺失突变,以及单一位点多种核苷酸/氨基酸替换突变所占比例.一致率是选定区域全部克隆均为一致的核苷酸/氨基酸数目与最长的单一克隆核苷酸/氨基酸序列数目之比,表示该段区域核苷酸/氨基酸序列一致性,用于展示区域内部的替换突变或/和缺失突变.如插入突变的克隆不能占据简单多数,则阳性率和一致率均为0%. Vector 8.0版软件比较核苷酸/氨基酸可推导出系统发生树,以分支树形式表示分子之间的进化关系,以及遗传关系的远近.

2 结果

2.1 HBV基因组 经过在GenBank中按不同基因型进行搜寻,分别获得的基因型A型病毒株12株,B型16株,C型48株,D型13株,E型4株,F型26株,G型12株,H型3株,共134株HBV病毒株全基因组序列.按照方法中的排除条件,剔除部分不适合本研究的病

表1 主要HBV长度类型的基因型分布序列一致性与阳性率

HBV 不同基因组	3215 nt	3182 nt	长度范围(nt)	阳性率(%)	一致率(%)
A型(13株)	0	0	3 149-3 254	98.8	85.3
B型(11株)	9	0	3 194-3 221	99.7	88.5
C型(41株)	31	4	2 996-3 215	99.7	70.4
D型(16株)	0	12	3 182-3 194	99.5	86.4
E型(3株)	0	0	3 212	99.9	96.8
F型(22株)	18	0	3 129-3 215	100.0	84.0
G型(9株)	1	0	3 089-3 248	100.0	93.2
H型(4株)	3	0	3 206-3 215	100.0	95.9

毒株序列, 共有119株HBV病毒株全基因组核苷酸序列引入到本研究中, 选取率为88.8%。选用的序列中基因型A型12株, B型11株, C型44株, D型13株, E型3株, F型21株, G型12株, H型3株。其中C型的44株病毒株包括本组以往的研究^[15-16]获得的5个HBV全基因组克隆, 分别命名为China Dong 1 C, 2 C, 3 C, 6 C和7 C。在本文中选择的病毒株的编号为地点+顺序号+基因型别, 如Japan 12 C, Japan是该序列来源国家, 12为本研究组在GenBank搜寻过程中所定义的顺序号, C为基因型别, 按照此方法, 所有参加比较的序列各有其独立的标记。本研究组既往研究中获得的序列在国家后加Dong以示区别。基因型别按照报告单位的地理分布如下: A型: 加拿大1株, 法国4株, 南非7株; B型: 中国3株, 日本2株, 荷兰1株, 南非1株, 瑞典4株; C型: 澳大利亚5株, 中国17株, 日本13株, 瑞典9株; D型: 西班牙1株, 法国2株, 日本5株, 瑞典5株; E型: 加纳1株, 日本2株; F型: 阿根廷4株, 瑞典6株, 委内瑞拉11株; G型: 德国1株, 日本11株; H型: 瑞典3株。

2.2 核苷酸序列的一致性 HBV基因组多态性表现为长度的明显差异, 所研究的119株病毒株中, 长度最长为3 254 nt, 最短为2 996 nt; 62株病毒株的基因组全长为3 215 nt, 占研究总数的52.1%(表1)。16株病毒株的基因组全长为3 182 nt, 占13.5%。HBV基因组长度具有一定的基因型特异性, 在GenBank中搜寻获得的12株G型HBV基因组序列中, 有7株病毒株的长度为3 248 nt, 该长度是G基因型的一个重要特征; 长度为3 221 nt的6个病毒株中, A型占5株; 长度为3 212 nt的3个病毒株均为E型。所研究的119株病毒株中有22种长度形式, 其中12种长度形式仅有1株病毒株(10.1%, 12/119)。

所有119株HBV基因组序列比较后, 推导出的系统发生树。经过比较后, 可以得出以下几点结论: (1)除C基因型病毒株之外, 各基因型的HBV均表现为相对独立的分支, C型在各基因型中显得较为古老; (2)虽然

在存储HBV基因组序列时, 各作者进行了基因型的判断, 但经过全基因组序列分析, 发现部分分型方法并不正确, 如Australian 3 C, 4 C和5 C应当属于D基因型, Japan 1 G属于F基因型, Japan 11 G属于H基因型, Japan 2 G属于A基因型。因此将上述型别重新划分到各基因型组中, 进行比较分析。

在119个基因型序列的HBV基因组比较后, 基因组的跨度总长为3 303 nt, 这是由于不同型别之间在不同区域存在插入突变, 导致基因组序列长于最长的3 254 nt。比较后总阳性率为95.7%, 说明有140 nt (4.3%, 140/3 303)为插入或缺失突变; 总一致性仅为47.7%(表2), 说明一半以上的位点存在多态性表现。

表2 119例不同基因型HBV病毒株基因组序列一致性比较

	阳性率(%)	一致率(%)
1-700	88.1	30.6
701-1 102	96.3	53.1
1 103-1 653	97.1	20.8
1 654-1 950	99.3	74.2
1 951-3 303	98.1	59.0
总计	95.7	47.7

2.3 蛋白质序列的一致性分析 应用Vector 8.0版软件ORF判定功能, 对HBV基因组蛋白编码区域进行分析, 发现40株序列编码前前-S多肽, 占33.6%, 分别来自C, F和H基因型(Japan11G划归H基因型, 下同)。17株序列编码前-X多肽, 占14.3%, 均来自C基因型, 4株来自日本病毒株, 13株来自中国病毒株; 其中16株病毒序列编码前前-S多肽, 占94.1%(16/17)。31株病毒序列不编码前-C多肽, 占26.1%, 除了E, F基因型外, 各基因型均有病毒株不编码前-C多肽。根据HBV基因组编码病毒蛋白结构的不同, 将HBV分为7种蛋白表型, 分型方式见表4。根据这种分型方式, 我们分析的119株病毒基因组中, I型: 共10株, 均来自C基因型, 占8.40%; II型: China Dong 6C, 7C, China

14C, 共3株, 占2.5%; III型: China 1C, 仅1株, 占0.8%; IV型: 共47株, 占39.5%, 为主要的流行蛋白表型, 来自除G, H基因型以外的所有基因型; V型, 23株, 占19.3%, 来自C, D, F和H基因型; VI型: Australian 2C, Japan 4C, Sweden 1C, Japan 11G(H型), 共4株, 占3.4%; VII型: 23株, 占19.3%, 来自A, B, C, D和G基因型, 其中G基因型前-C多肽的表达具有其独特的特异性, 长度小于传统定义的前-C区29 aa长度的多肽, 仅长12 aa; 有8个病毒株不能完整表达一种或一种以上HBV病毒蛋白, 表达序列长度小于预计长度的50%, 占6.7%, 无法进行蛋白表型分型, 属于缺陷型病毒。我们提出的HBV蛋白表型分型方法, 以IV型为主要流行表型, 占39.5%, 为Galibert *et al*^[2]最早解读的HBV基因组序列编码方式; V型和VII型各占19.3%, 是重要的流行型别。前-X多肽多与前-S多肽联动表达, 仅1例例外。I, II, III型均来自C基因型, I型以China Dong 1C为代表, 占8.4%。

病毒蛋白的表达除部分位点上表现出型特异性特征外, 也表现出区域性的多态性, 113株不同基因型HBV编码的HBeAg总阳性率为98.6%, 总一致率仅为37.4%, 其大部分区域存在明显的氨基酸位点多态性(见表3)。分区的比较发现: 前-C区编码的29 aa, 31个序列不编码前-C多肽, 故其一致率为0%; 30-61aa和130-173aa为HBeAg的2个高保守区, 其阳性率均为100%, 无明显的插入/缺失突变; 这2段区域的一致率明显高于总一致率。相应的, 62-129 aa和174-214 aa为HBeAg的2段高变区, 一致率较总一致率分别低10.9%和10.3%。118株不同基因型HBV编码的全S基因编码产物包含前前-S区, 前-S1, 前-S2和主蛋白(下同), 总阳性率仅为87.3%, 总一致率仅为24.1%, 说明存在较多的缺失/插入突变(表3)。进一步分析发现存在2个高度变异区和1个高度保守区, 以往定义的前-S1和前-S2

区为高变区, 46-187 aa区域的阳性率为95.8%, 一致率仅为7.7%; 410-448 aa一致率为0, 说明2个区域的缺失/插入突变较多。

118株不同基因型HBV编码的全X蛋白包括: 前-X和原X蛋白, 总阳性率仅为57.2%, 总一致率仅为27.7%, 说明存在较多的缺失/插入突变。进一步分析发现存在1个高度变异区和1个高度保守区。17株C基因型序列编码前-X多肽, 该段的阳性率和一致率为0%; 57-133 aa区域的阳性率为97.4%, 一致率低为54.5%, 高于总一致率1倍; 134-264 aa一致率为23.7%, 说明区域内存在较多的缺失/插入突变。117株不同基因型HBV编码的多聚酶总阳性率仅为95.2%, 总一致率为43.5%, 其一致率为4个病毒蛋白中最高的, 与其编码区域内不包括表达前导区有关。进一步分析发现存在2个高度变异区和1个高度保守区, 180-369 aa区域的阳性率为87.4%, 一致率低为12.6%; 465-514 aa一致率为20.0%, 说明2个区域的缺失/插入突变较多。而370-464 aa区域的阳性率为99.0%, 一致率为76.8%, 远高于平均的43.5%。上述4种病毒蛋白按基因型进行比较后, 获得的阳性率和一致率见表4。将不同基因型病毒株编码的病毒蛋白进行比较, 所获得的系统发生树见图1A-E。

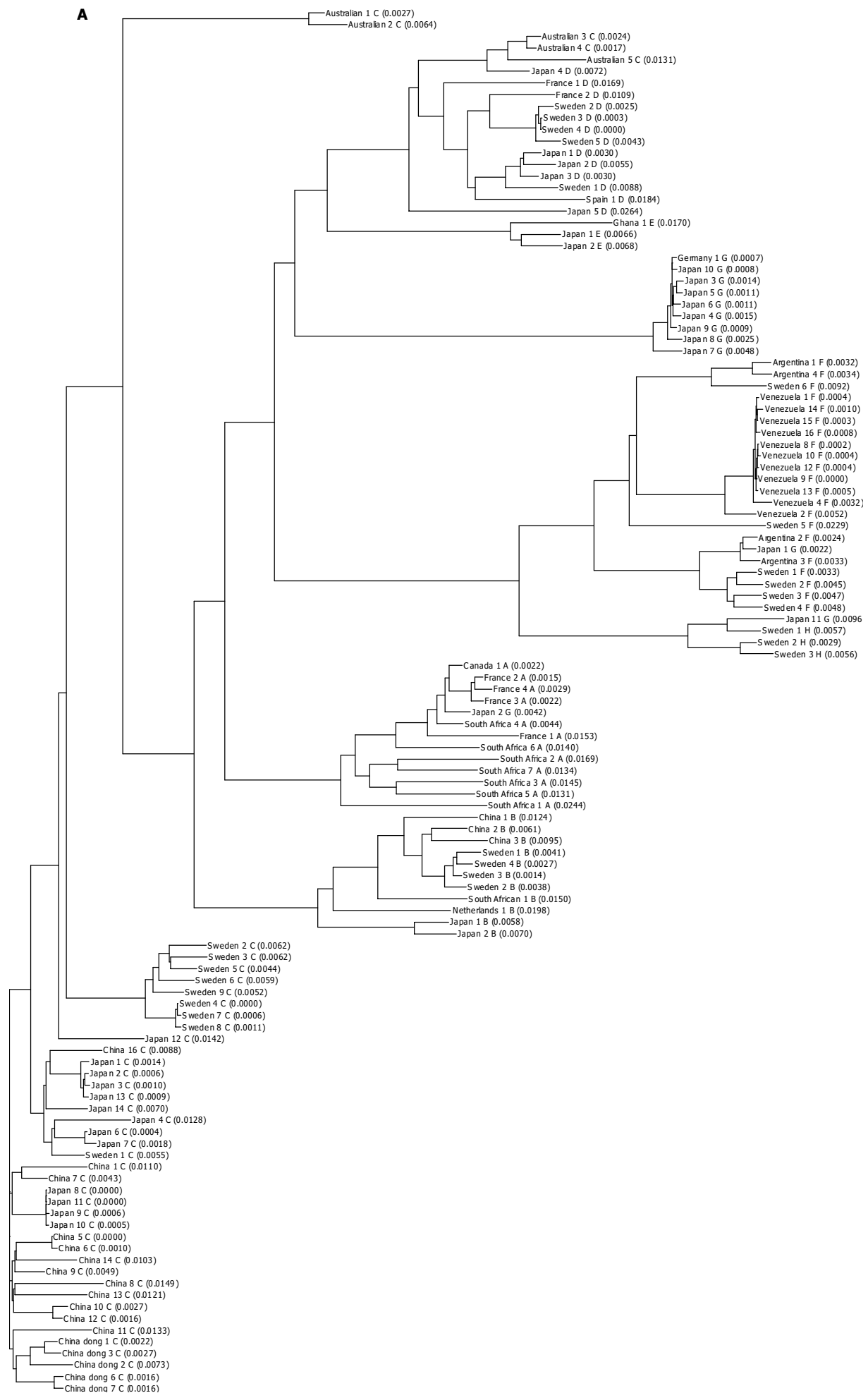
表3 113株HBV编码HBeAg和118株编码全S蛋白多态性的比较(%)

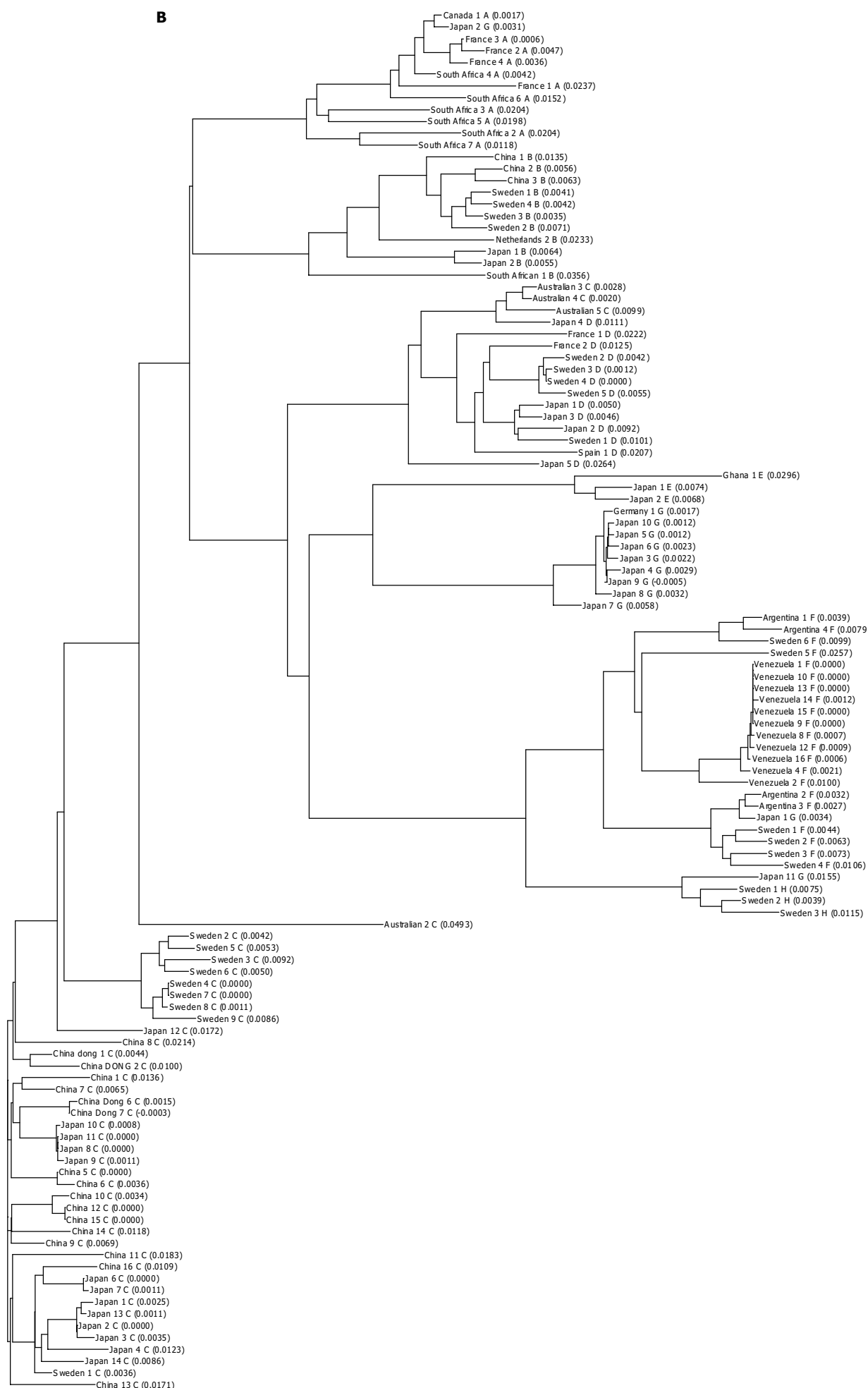
	HbeAg		全S蛋白	
	阳性率	一致率	阳性率	一致率
1-29	100.0	0.0	1-45	0.0
30-61	100.0	65.6	46-187	95.8
62-129	98.5	26.5	188-409	97.7
130-173	100.0	70.5	410-448	97.4
174-214	95.1	24.4		
总计	98.6	37.4	总计	87.3

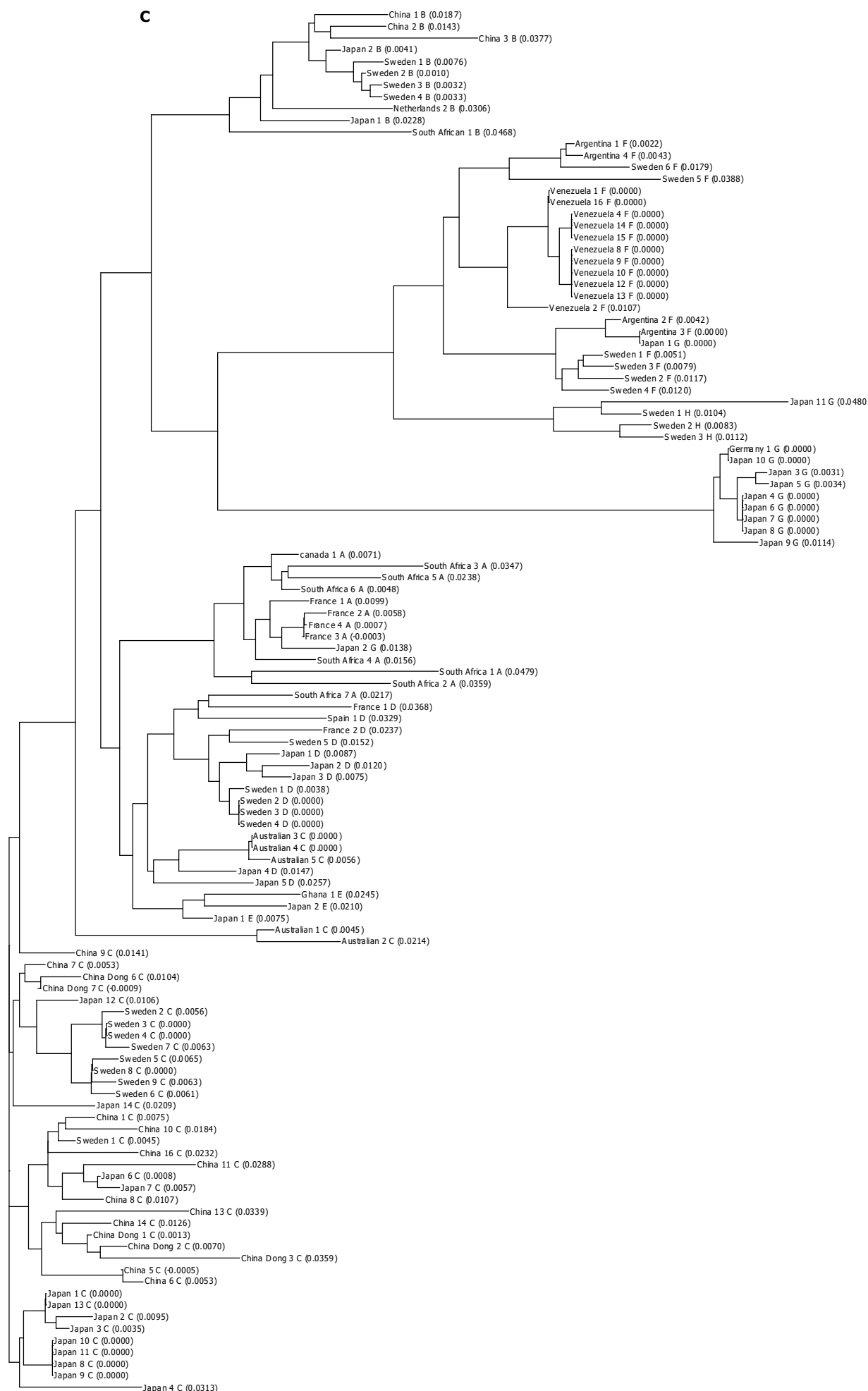
表4 HBV不同基因型病毒蛋白序列一致性与阳性率¹

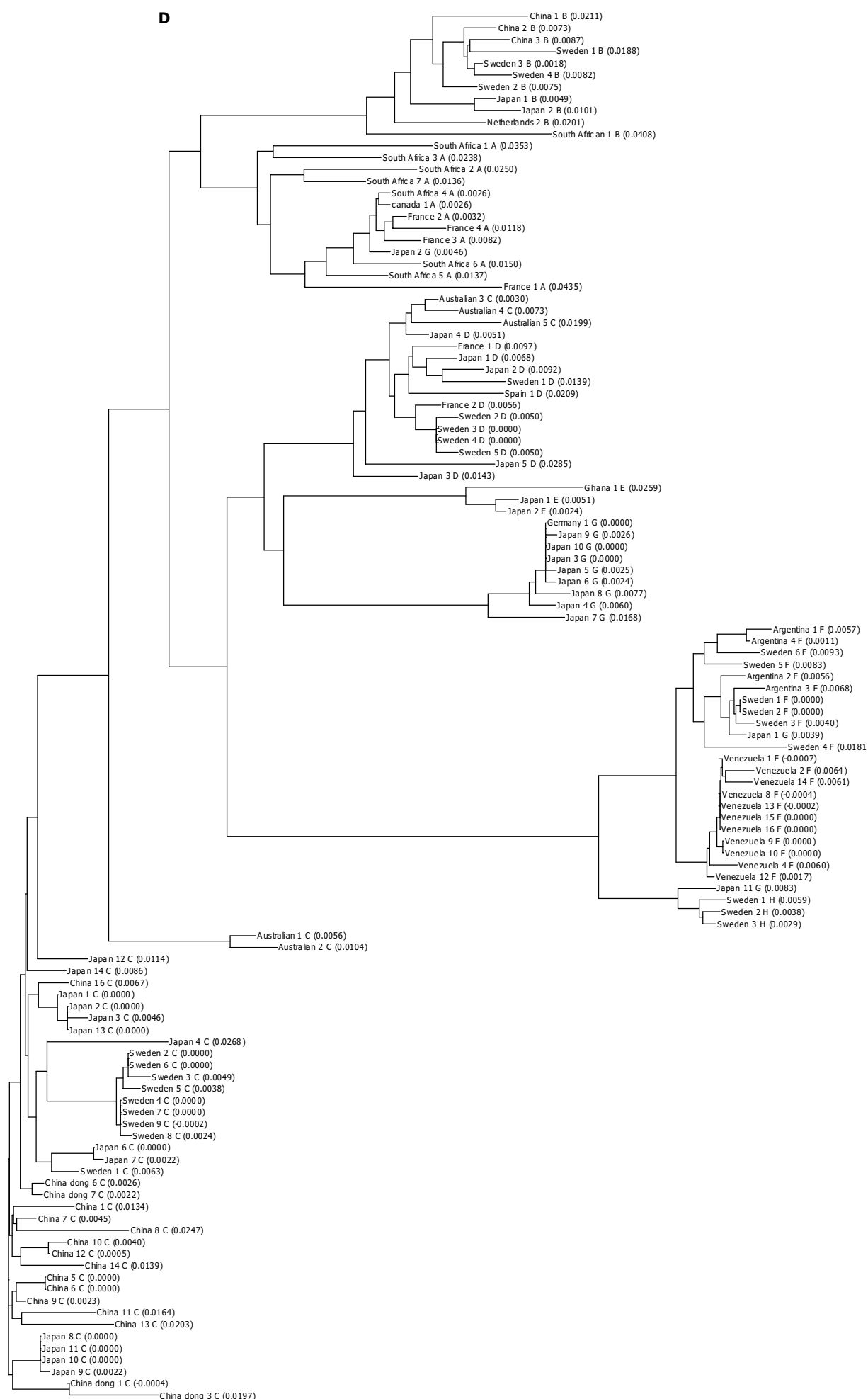
	全S蛋白		全X蛋白		全前C-C蛋白		多聚酶	
	阳性率(%)	一致率(%)	阳性率(%)	一致率(%)	阳性率(%)	一致率(%)	阳性率(%)	一致率(%)
A型(13株)	99.8	75.6	74.8	61.7	99.5	44.7	99.8	58.4
B型(11株)	100.0	81.5	96.8	82.7	99.5	77.4	99.8	85.5
C型(41株)	99.6	49.9	65.8	47.0	100.0	67.5	99.8	71.6
D型(16株)	98.7	85.3	100.0	83.8	99.5	69.8	99.4	85.5
E型(3株)	100.0	96.0	100.0	94.2	100.0	98.6	100.0	95.2
F型(22株)	89.9	80.0	100.0	82.5	100.0	76.9	100.0	84.9
G型(9株)	100.0	87.7	100.0	96.8	100.0	85.1	100.0	91.9
H型(4株)	100.0	96.9	100.0	89.6	100.0	84.0	100.0	95.1
总计	87.3	24.1	57.2	27.7	98.6	37.4	95.2	43.5

¹ 本表中所列的基因型为本文调整后的基因型。









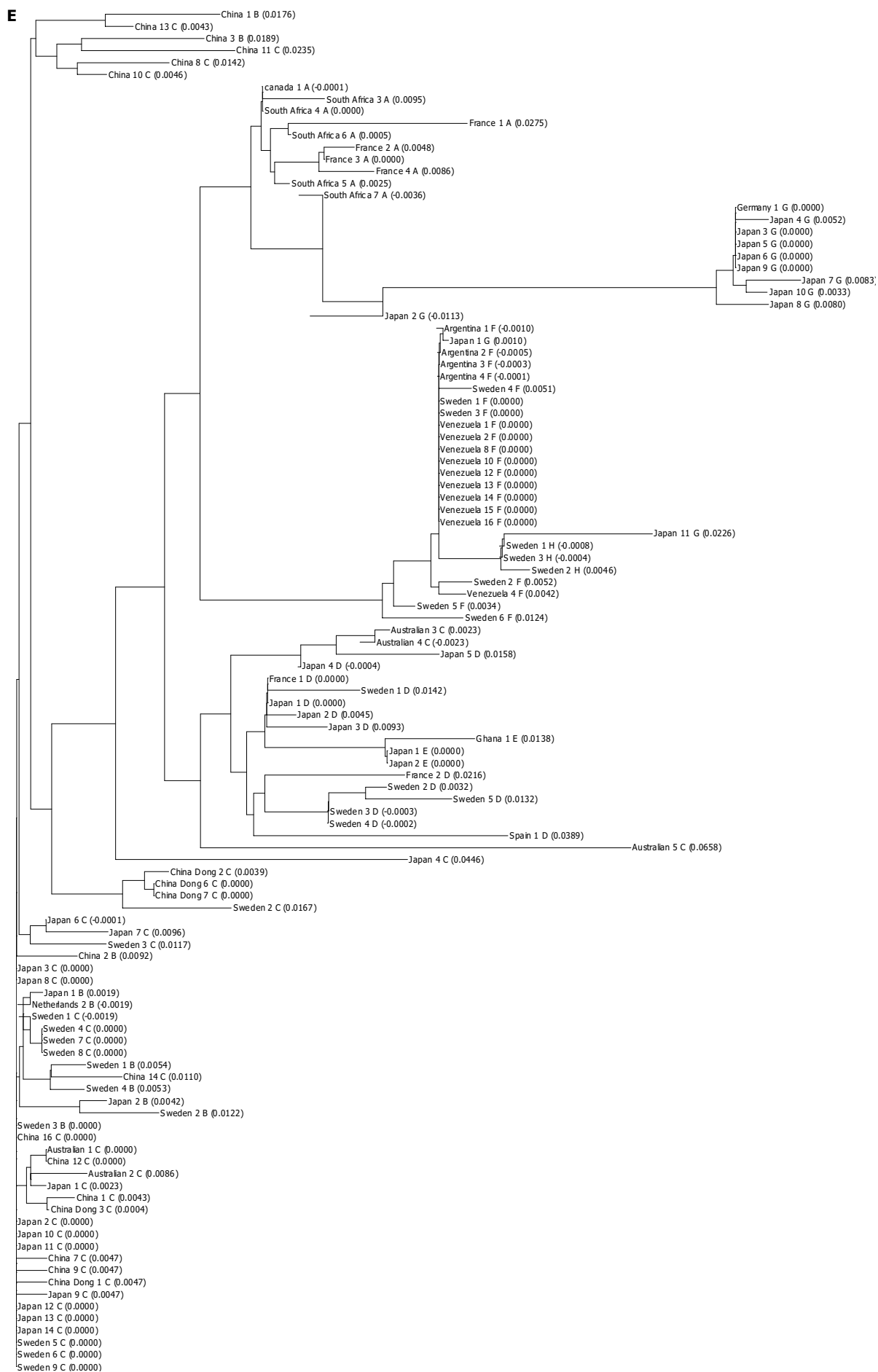


图1 A: HBV 基因型核苷酸序列系统发生树; B: 多聚酶氨基酸序列系统发生树; C: X 蛋白氨基酸序列系统发生树; D: HBV 表面蛋白氨基酸序列系统发生树; E: HBeAg 氨基酸序列系统发生树。

3 讨论

HBV 的基因组是最早被解读的病原体基因完整信息,我们探讨过 HBV 准种现象^[22-24],观察了 HBsAg 基因多态性与蛋白质多态性之间的关系^[25],HBV 血清型仅反映表面抗原的变化,不足以反映高度变异的病毒准种群核酸序列差异与病毒蛋白氨基酸序列差异之间的关系,现试图探讨一种以病毒蛋白差异为主要分型标准的 HBV 分型方式,暂命名为 HBV 蛋白表型(prototype).在 GenBank 中,我们初步搜寻出 200 多株 HBV 全基因组序列,之后再次限定作者标定的 HBV 基因型,自 A 至 H 进行再次筛选,选择出的病毒株按照不同的基因型分组以备进一步分析.将筛选出的 HBV 基因组序列进行初步分析,凡是序列内部包含有测序结果不精确的病毒株均被排除在本研究之外.经过上述选择过程,共筛选出 119 株 HBV 全基因组序列,长度最长为 3 254 nt,最短为 2 996 nt. HBV 多态性表现在长度的一致性,共有 22 种长度形式,52.1% 的病毒株基因组长度为 3 215 nt,为主要流行长度;13.5% 的病毒株基因组全长为 3 182 nt;其中 12 种长度形式仅有 1 株病毒株.基因组长度具有一定的基因型特异性,经过调整后, G 基因型 9 例病毒株中, 7 例长度为 3 248 nt, 经比较发现在 91 nt 之后有一段长 26 nt 的插入序列, 为 5' TAGAACAACCTTTGCCATATGGCCTTTTGGCTTAGA-3' G 基因型特异性序列. 该段序列与上游 5 nt 共同编码前 C 区 MDRTTLPYGLFGL, 成为缩水的前 -C 区, 替代了其他基因型中的前 -C 区编码的 29 aa 和 HBcAg 的第一位 M. 该编码现象是 G 基因型的一个重要蛋白分子特征.

基因型的分型原则是根据全基因核苷酸序列之间比较后所获得的一致性进行分类的,至此,目前报告了 8 种基因型^[3, 4, 6, 7, 26]. 我们自 GenBank 中搜索出的 119 个作者进行基因型分型的基因组病毒株中,调整后 A 型 13 株, B 型 11 株, C 型 41 株, D 型 16 株, E 型 3 株, F 型 22 株, G 型 9 株, H 型 4 株,其中以 C 型最多,占 34.5%;其次为 F 型,占 18.5%. 曾有学者^[7, 26]认为基因型的分布具有明显的地理特性,但本研究 119 株病毒株的地理分布广泛,欧洲 36 株,亚洲 53 株,非洲 9 株,美洲 16 株,大洋洲 5 株.除报告较少的 E 型和 H 型, HBV 基因型的地理分布没有一定的规律,可能与目前世界范围交通发达,人员流动大有关,但在与外界交流很少的民族中可能存在少见的 HBV 基因型. 我们针对 119 例不同基因型的 HBV 基因组进行了全序列比较,同时对其编码的不同的病毒蛋白进行了比较分析,结果发现: (1)除全 C 蛋白氨基酸系统发生树外,基因组系统发生树与全 S 蛋白、多聚酶和全 X 蛋白的氨基酸系统发生树结构相似,而在全 C 蛋白的氨基酸系统发生树分析过程中,发现部分 B 基因型病毒株与 C 基因型的全 C 蛋白的系统发生特征相近. 从上述图形结构的比较而言,可以认为限定区域

内核核苷酸的变异导致的氨基酸序列的变异程度是不均一的,两种分子的进化步骤具有非同步性; (3)无论基因组核苷酸系统发生树还是 4 种病毒蛋白的系统发生树,各基因型的 HBV 病毒株均划归于相对独立的分支,表现出较高的遗传特征,唯有 C 型各病毒株的遗传关系显得较为松散,其系统发生树的地位在各基因型中显得较为古老; (3)结合系统发生树与基因型内病毒株序列一致性的研究结果发现: 只有 C 基因型各病毒株之间的总一致率低为 70.4%(表 1),其他基因型的总一致率接近或高于 84.0%. 研究同时发现随着分型内部病毒株的数量增加,其一致率出现明显下降, C 型一致率的水平已低于以往定义的基因型的分型标准,按目前标准定义的 C 基因型的病毒株总体的遗传特征不明确,有必要提出新的分型标准来表示 HBV 的变异特征.

我们进一步分析了 HBV 不同基因型的节段差异性,发现 119 株不同基因型 HBV 病毒株基因组序列总一致率为 47.7%. 区域 1-700 nt 一致率 30.6%, 区域 1 103-1 653 nt 一致率 20.8%, 较总一致率低 10-20%, 为高变区; 区域 1 654-1 950 nt 的一致率为 74.2%, 可定为高保守区. 我们是在计算机软件分析的基础上,人工的最小化或最大化各区域的一致率,具体流程见文献^[18],但划分出来的高变区和高保守区与文献^[18]不同,这可能与选择的基因型与血清型不同有关. 本文还进一步确定了各病毒蛋白的高变区和高保守区, 113 株全 C 蛋白多态性的比较结果提示 62-129 aa 和 174-214 aa 一致率分别为 26.5% 和 24.4%; 30-61 aa 和 130-173 aa 一致率分别为 65.6% 和 70.5%, 较总一致率 37.4% 有明显区别. 118 例全 S 蛋白内多态性的比较结果提示 46-187 aa 和 410-448 aa 一致率分别为 7.7% 和 0.0%; 188-409 aa 一致率 43.7%, 较总一致率 24.1% 有明显区别. 118 例全 X 蛋白多态性的比较结果提示 57-133 aa 一致率 54.5%, 明显高于总一致率的 27.7%; 117 例多聚酶多态性的比较结果提示 180-369 aa 和 465-514 aa 一致率分别为 12.6% 和 20.0%; 370-464 aa 一致率 76.8 %, 较总一致率 43.5% 有明显区别, 这些数据说明各病毒蛋白内部均有各自特异的高变区和高保守区. 全 S 蛋白总一致率仅 24.1%, 其前 S1、前 S2 区一致率小于 10%, 目前疫苗的靶区域 188-409 aa 一致率也仅 43.7%, 这提示目前疫苗的覆盖性较差, 需要进一步研制高代表性 HBV 疫苗.

我们着重研究了 HBV 4 个 ORF 的结构特征,并提出根据病毒蛋白不同的结构特征进行分型的新概念,目前暂时将这种型别命名为 HBV 的蛋白表型. 在分型方法的建立过程中,我们着重强调了前前 -S 多肽、前 -X 多肽和前 -C 多肽的重要性,这是由于这 3 段蛋白序列的编码与否在不同病毒株基因组的表现形式是不一样的,前前 -S 多肽长度为 45 aa, 占全 S 蛋白全长的 10.1%(45/446), 编码前前 -S 多肽阳性的病毒株占研究总数的 33.6%, 分别属于 C, F 和 H 基因型; 前 -X 多

肽长度为 56 aa, 占全 X 蛋白全长的 26.7%(56/210), 编码前 -X 多肽阳性的病毒株占研究总数的 14.3%, 均属于 C 基因型;前 -C 多肽长度 29 aa, 占前 C-C 蛋白全长的 13.7%(29/212), 编码前 -C 多肽阴性的病毒株占研究总数的 26.1%, 分别属于 A, B, C, D, F 和 H 基因型, G 型的 9 株病毒株的前 -C 区编码多肽形式特别. 我们发现编码前前 -S 多肽和前 -X 多肽在病毒基因组中属于相对少见的情形, 而编码前 -C 多肽在病毒基因组中属于较常见的现象. 我们以往关于前 -C 区与 C 区基因相互关系的研究提出前 -C 区编码与 HBeAg 的生物合成有关, 前 -C 区变异^[29]或核心蛋白启动子(CP)的变异^[30]可导致 HBeAg 阴性慢性乙型肝炎(CHB); 前 -X 多肽的编码与原发性肝癌(HCC)的发生有关^[31], 但前前 -S 多肽的功能尚不明. 蛋白表型的分型方法是确定的, 这与基因型分型日益增多的情形大相径庭, 由于基因型表现出明显的地理特性, 随着研究的深入, 其分型结果会日益增多, 由于蛋白表型的分型方式固定, 将简化分型方式, 为以后的研究设立一个便于比较的平台.

蛋白表型分型重点强调由于基因突变所导致原定义的 HBV 表面抗原、X 蛋白、核心蛋白的前导性序列以及病毒蛋白结构的变化, 同时也强调了病毒蛋白前导性序列的重要性. 按照以往的概念, HBV 表面抗原、核心蛋白为病毒的结构蛋白, 其前导性序列可能与病毒蛋白合成后的细胞定位有关; X 蛋白的功能尚不明, 可能与病毒基因调节有关, 作为一种反式激活因子, 定位是在细胞核内. 蛋白表型分型强调了基于 HBV 病毒蛋白基本情况之上的变异方式, 显示核苷酸变异对病毒蛋白编码所产生的重要影响. 117 株病毒株的多聚酶的总一致率为 43.5%, 高于全 S 蛋白、全 X 蛋白和前 C-C 蛋白的总一致率 24.1%, 27.7% 和 37.4%, 这证实前导性序列对病毒蛋白一致性的影响. 按照我们提议的新的分型方法分析本研究搜集的 119 株病毒株, 发现为 Galibert *et al*^[12]最早解读的 HBV 基因组序列, 即 IV 型主要流行表型, 共 47 例, 占分析总数的 39.5%. IV 型的蛋白表型特点是表达前 -C 多肽, 而不表达前前 -S 多肽和前 -X 多肽, 除 G, H 基因型外, 其他基因型病毒株均可见 IV 型的分布. V 型和 VII 型各有 23 例, 占分析总数的 19.3%, 是重要的流行型别, 前者的特点是表达前前 -S 多肽和前 -C 多肽, 但不表达前 -X 多肽, 后者的特点是三种前导性序列均不表达. II 型, III 型和 VI 型是少见的蛋白表型, 其中 III 型仅 1 例, 其存在与否需要进一步证实. 由于 HBV 基因存在高变异现象, 1993 年有学者提出 HBV 准种学说^[32-33], 我们的研究^[15]证实 HBV 在患者体内以准种群形式存在, 这种学说同时带来一种研究技术上的要求, 即以单独克隆的 HBV 基因组测序结果为准, 而不应当是以往研究^[34-36]中以直接的 PCR 结果进行测序来表示全基因组情况, 因为 PCR- 直接测序方法不能排除准种群

导致的位点混杂情况. 有鉴于此, II 型、III 型和 VI 型的病毒株测序方法存有一定误差, 其可靠性需要进一步验证, 因此按照本文提出的分型方法, 分型的型别可能将有所减少. 新的分型结果发现有 8 株 HBV 基因组序列不能编码完整的 1 种或 1 种以上的病毒蛋白, 推断其不具有独立完成生活史的能力, 即需要依靠其他 HBV 病毒株生存的缺陷型病毒, 我们^[37]以往的研究也发现这种现象, 按照本方法无法进行分型, 因此设立无法分型的型别为缺陷型病毒特异性型别. 49 株亚洲 HBV 病毒分离株的蛋白表型分型结果提示: I, IV, V 和 VII 型 HBV 病毒株所占比例均大于 20%, 7 种型别均出现于亚洲; 36 株欧洲 HBV 病毒分离株的蛋白表型分型结果提示: IV 型占 58.3%, VII 型占 25.0%, V 型占 13.9%, 说明以 IV 型为主要流行型别, VII 型也占重要地位. I, II, III 型不出现在欧洲, 在亚洲这三种类型的总和占总数的 28.6%, 这 3 种类型的共同点是 HBV 基因组编码前 -X 多肽, 这 3 种蛋白表型 HBV 在亚洲的流行是否是亚洲原发性肝癌发病率高于欧洲的一个原因, 还需要进一步研究.

总之, 我们应用生物信息学技术对 GenBank 中存储的 119 个不同基因型 HBV 病毒株全基因组进行了比较分析, 在此基础上我们提出 HBV 蛋白表型假说, 蛋白表型着重强调了 S 基因、X 基因、C 基因前导性序列的重要性, 通过分型将局部核苷酸替换突变与蛋白结构差异结合起来, 我们试图建立一种区别于现有基因型分型方式, 以展示病毒蛋白结构差异为主的分型方式, 这种分型方式的现实意义需要进一步探讨.

4 参考文献

- 1 Le Bouvier GL, McCollum RW, Hierholzer WJ Jr, Irwin GR, Krugman S, Giles JP. Subtypes of Australia antigen and hepatitis-B virus. *JAMA* 1972;222:928-930
- 2 Galibert F, Mandart E, Fitoussi F, Tiollais P, Charnay P. Nucleotide sequence of the hepatitis B virus genome (subtype ayw) cloned in *E. coli*. *Nature* 1979;281:646-650
- 3 Okamoto H, Tsuda F, Sakugawa H, Sastrosoewignjo RI, Imai M, Miyakawa Y, Mayumi M. Typing hepatitis B virus by homology in nucleotide sequence: comparison of surface antigen subtypes. *J Gen Virol* 1988;69(Pt 10):2575-2583
- 4 Norder H, Hammas B, Magnius LO. Typing of hepatitis B virus genomes by a simplified polymerase chain reaction. *J Med Virol* 1990;31:215-221
- 5 Kao JH, Chen PJ, Lai MY, Chen DS. Hepatitis B genotypes correlate with clinical outcomes in patients with chronic hepatitis B. *Gastroenterology* 2000;118:554-559
- 6 Stuyver L, De Gendt S, Van Geyt C, Zoulim F, Fried M, Schinazi RF, Rossau R. A new genotype of hepatitis B virus: complete genome and phylogenetic relatedness. *J Gen Virol* 2000;81(Pt 1): 67-74
- 7 Arauz-Ruiz P, Norder H, Robertson BH, Magnius LO. Genotype H: a new Amerindian genotype of hepatitis B virus revealed in Central America. *J General Virol* 2002;83(Pt 8): 2059-2073
- 8 黄晶, 高志良. 乙型肝炎病毒基因型及其临床意义的研究. *世界华人消化杂志* 2002;10:1362-1364
- 9 温志立, 谭德明. 多对型特异性引物巢式 PCR 检测湖南省乙肝病毒基因型. *世界华人消化杂志* 2004;12:332-335
- 10 Sugauchi F, Orito E, Ichida T, Kato H, Sakugawa H, Kakumu S, Ishida T, Chutaputti A, Lai CL, Gish RG, Ueda R, Miyakawa

- Y, Mizokami M. Epidemiologic and virologic characteristics of hepatitis B virus genotype B having the recombination with genotype C. *Gastroenterology* 2003;124:925-932
- 11 董菁, 成军, 杨倩. 乙型肝炎病毒新开放读码框架的确定及其意义. 世界华人消化杂志 2004;12:757-762
 - 12 董菁, 成军. 乙型肝炎病毒前S区基因的界定. 世界华人消化杂志 2003;11:1091-1096
 - 13 董菁, 成军. 乙型肝炎病毒前-X基因的初步研究. 世界华人消化杂志 2003;11:1097-1101
 - 14 董菁, 洪源, 刘妍, 钟彦伟, 王琳, 王刚, 张玲霞, 陈菊梅. 乙型肝炎病毒中国流行株全基因的克隆与序列分析. 世界华人消化杂志 2003;11:1119-1126
 - 15 成军, 董菁. 乙型肝炎病毒基因组结构与功能复杂性的新认识. 世界华人消化杂志 2003;11:1073-1080
 - 16 董菁, 成军, 王勤环, 皇甫竞坤, 施双双, 张国庆, 洪源, 李莉, 斯崇文. 慢性乙型肝炎患者体内乙型肝炎病毒DNA序列异质性及准种特点的研究. 中华医学杂志 2002;82:81-85
 - 17 董菁, 成军, 皇甫竞坤, 洪源, 王刚, 陈国凤, 李莉, 张玲霞, 陈菊梅. 乙型肝炎病毒序列准种个体化特征的研究. 解放军医学杂志 2002;27:119-121
 - 18 Gunther S, Li BC, Miska S, Kruger DH, Meisel H, Will H. A novel method for efficient amplification of whole hepatitis B virus genomes permits rapid functional analysis and reveals deletion mutants in immunosuppressed patients. *J Virol* 1995;69:5437-5444
 - 19 董菁, 成军, 杨倩, 纪冬, 张健, 李莉. 乙型肝炎病毒基因组高变区界定的初步研究. 世界华人消化杂志 2004;12:42-46
 - 20 杨倩, 董菁, 成军, 刘妍, 洪源, 王建军, 王琳, 张树林. 乙型肝炎病毒基因组中前-X-编码基因启动子序列的确定及转录活性的鉴定. 解放军医学杂志 2003;28: 763-765
 - 21 杨倩, 董菁, 成军, 刘妍, 洪源, 王建军, 张树林. 乙型肝炎病毒基因组中前-X-编码基因启动子序列的确定及转录活性的鉴定. 解放军医学杂志 2003;28:761-762
 - 22 董菁, 李进, 施双双, 皇甫竞坤, 成军, 王勤环, 洪源, 李莉. 乙型肝炎病毒基因组准种与变异特点的研究. 解放军医学杂志 2002;27:116-118
 - 23 董菁, 施双双, 皇甫竞坤, 成军, 王勤环, 李莉, 斯崇文. 乙型肝炎病毒X基因准种特点的研究. 中国病毒学 2002;17:22-26
 - 24 董菁, 成军, 王勤环, 施双双, 洪源, 皇甫竞坤, 王刚, 李莉, 斯崇文. 乙型肝炎病毒逆转录酶区基因序列准种与变异研究. 解放军医学杂志 2001;26:823-825
 - 25 董菁, 刘妍, 皇甫竞坤, 施双双, 王刚, 洪源, 陈国凤, 李莉, 陈菊梅, 成军. 乙型肝炎病毒表面抗原一级结构多态性的初步研究. 胃肠病学和肝病杂志 2002;11:130-135
 - 26 Norder H, Hammas B, Lee SD, Bile K, Couroucé AM, Mushahwar IK, Magnus LO. Genetic relatedness of hepatitis B viral strains of diverse geographical origin and natural variations in the primary structure of the surface antigen. *J General Virol* 1993;74(Pt 7):1341-1348
 - 27 杨倩, 董菁, 成军. 乙型肝炎病毒前-S基因的分子流行病学研究. 世界华人消化杂志 2004;12:785-789
 - 28 董菁, 杨倩, 成军. 乙型肝炎病毒前-X基因的分子流行病学研究. 世界华人消化杂志 2004;12:794-800
 - 29 董菁, 施双双, 皇甫竞坤, 成军, 王勤环, 王刚, 洪源, 李莉, 斯崇文. 乙型肝炎病毒前C/C基因准种与变异特点的研究. 解放军医学杂志 2002;27:122-124
 - 30 董菁, 施双双, 张国庆, 皇甫竞坤, 洪源, 成军, 王勤环, 李莉, 斯崇文. 乙型肝炎病毒C基因启动子区异质性检测初步研究. 临床检验杂志 2002;20:72-74
 - 31 Takahashi K, Akahane Y, Hino K, Ohta Y, Mishiro S. Hepatitis B virus genomic sequence in the circulation of hepatocellular carcinoma patients: comparative analysis of 40 full-length isolates. *Arch Virol* 1998;143:2313-2326
 - 32 Blum HE. Hepatitis B virus: significance of naturally occurring mutants. *Intervirology* 1993;35:40-50
 - 33 Carman W, Thomas H, Domingo E. Viral genetic variation: hepatitis B virus as a clinical example. *Lancet* 1993;341:349-353
 - 34 Hannoun C, Horal P, Lindh M. Long-term mutation rates in the hepatitis B virus genome. *J Gen Virol* 2000;81:75-83
 - 35 Kramvis A, Weitzmann L, Owiredo WK, Kew MC. Analysis of the complete genome of subgroup A' hepatitis B virus isolates from South Africa. *J Gen Virol* 2002;83(Pt 4):835-839
 - 36 Sugauchi F, Mizokami M, Orito E, Ohno T, Kato H, Suzuki S, Kimura Y, Ueda R, Butterworth LA, Cooksley WG. A novel variant genotype C of hepatitis B virus identified in isolates from Australian Aborigines: complete genome sequence and phylogenetic relatedness. *J Gen Virol* 2001;82(Pt 4):883-892
 - 37 董菁, 成军, 王勤环, 王刚, 施双双, 夏小兵, 斯崇文. 外周血中乙型肝炎病毒截短型囊膜蛋白基因的克隆化与序列分析. 中华肝脏病杂志 2001;9:163-165

ISSN 1009-3079 CN 14-1260/R 2004 年版权归世界胃肠病学杂志社

• 消息 •

《中国生物学文摘》收录WJG和世界华人消化杂志

本刊讯 经专家评估和遴选, *World Journal of Gastroenterology*(WJG)和世界华人消化杂志被《中国生物学文摘》和中国生物学文献数据库收录. 中国生物学文献数据库在期刊的基础上开发建设, 数据量已达20万多条, 并形成了期刊、光盘、网络版系列产品.《中国生物学文摘》1998年获得第六次全国科技期刊文献检索出版物评比一等奖.(世界胃肠病学杂志 2004-05-05)

Nature Clinical Practice Gastroenterology & Hepatology 收录 *World Journal of Gastroenterology*

本刊讯 在2004年11月, Nature Publishing Group 将会出版一系列杂志, 题为自然临床实践, 包括肿瘤、心血管、泌尿、胃肠病学和肝脏病学, 这些文章会应用于临床患者和医生. Nature Publishing Group 收录非常有影响的杂志, *World Journal of Gastroenterology* 也被收录. *Nature Clinical Practice Gastroenterology & Hepatology* 为非常繁忙的胃肠病学家和肝脏病学家提供其专业的概况和其领域的所有的关键的进展, 更重要的是这些进展会为他们患者提供进一步的帮助. *Nature Clinical Practice Gastroenterology & Hepatology* 提供电子和印刷版, 其主编为 Stephen Hanauer. (世界胃肠病学杂志 2004-06-15)