

DOI:10.12113/202001007

2019 新型冠状病毒基因组的生物信息学分析

陈嘉源¹, 施劲松², 丘栋安³, 刘畅⁴, 李鑫¹, 赵强¹, 阮吉寿⁵, 高山^{1*}

(1.南开大学 生命科学学院, 天津 300071; 2.东部战区总医院, 南京 210016;

3.英国诺丁汉特伦特大学 生物科学系, 诺丁汉 NG11 8NS;

4.南开大学 医学院, 天津 300071; 5.南开大学 数学科学学院, 天津 300071)

摘要:2019年12月,中国武汉报道了冠状病毒引起的肺炎,其临床症状与2003年爆发的严重急性呼吸综合征(Severe Acute Respiratory Syndrome, SARS)不同,因此推断该病毒可能是冠状病毒的一个新变种。不同于简单使用全基因组序列的其它研究,我们于2018年在国际上首次提出分子功能与进化分析相结合的研究思想,并应用于Beta冠状病毒B亚群(BB冠状病毒)基因组的研究。在这一思想指导下,本研究使用BB冠状病毒基因组中的一个互补回文序列(命名为Nankai complemented palindrome)与其所在的编码区(命名为Nankai CDS)对新发布的2019新型冠状病毒基因组(GenBank: MN908947)进行分析以期准确溯源,并对BB冠状病毒的跨物种传播和宿主适应性进行初步研究。溯源分析的结果支持2019新型冠状病毒源自蝙蝠,但与SARS冠状病毒差异巨大,这一结果与两者临床症状差异一致。本研究的最重要发现是BB冠状病毒存在大量的可变翻译,从分子水平揭示了BB冠状病毒变异快、多样性高的特点。从BB冠状病毒可变翻译中获取的信息可应用于(但不限于)其快速检测、基因分型、疫苗开发以及药物设计。另外,我们推断BB冠状病毒可能通过可变翻译以适应不同宿主。基于大量基因组数据的实证分析,本研究在国际上首次从分子水平尝试解释了BB冠状病毒变异快、宿主多且具有较强的宿主适应性的原因。

关键词:冠状病毒;SARS;可变翻译;SARS-CoV-2;跨物种传播

中图分类号:Q93 **文献标志码:**A **文章编号:**1672-5565(2020)02-096-07

Bioinformatics analysis of the 2019 novel coronavirus genome

CHEN Jiayuan¹, SHI Jinsong², YAU Tungon³, LIU Chang⁴

LI Xin¹, ZHAO Qiang¹, RUAN Jishou⁵, GAO Shan^{1*}

(1. College of Life Sciences, Nankai University, Tianjin 300071, China;

2. National Clinical Research Center of Kidney Disease, Jinling Hospital, Nanjing University School of Medicine, Nanjing 210016, China;

3. School of Science and Technology, Nottingham Trent University, Nottingham, NG11 8NS, United Kingdom;

4. School of Medicine, Nankai University, Tianjin 300071, China;

5. School of Mathematical Sciences, Nankai University, Tianjin 300071, China)

Abstract: In December 2019, a pneumonia outbreak caused by a human coronavirus was reported in Wuhan (China). This virus was predicted as a new coronavirus, named the 2019 novel coronavirus (2019-nCoV), as it caused clinical symptoms different from Severe Acute Respiratory Syndrome (SARS) during the 2003 outbreak. Currently, most of the researchers simply use the complete genome or specific structural gene sequences to investigate coronavirus (e.g. phylogenetic analysis) without considering the functions of the products from coronavirus genes. To overcome this shortcoming, we proposed the joint analysis of the molecular function and phylogeny, and applied it in our previous study of genomes of Betacoronavirus subgroup B (BB coronavirus). In

收稿日期:2020-01-14;修回日期:2020-01-20.

基金项目:天津市重点研发计划科技支撑重点项目(南开大学, No.19YFZCZY00500).

作者简介:陈嘉源,男,在读本科生,研究方向:生物信息学.E-mail: 1610752@mail.nankai.edu.cn.

* 通信作者:高山,男,副教授、硕导,研究方向:生物信息学.E-mail: gao_shan@mail.nankai.edu.cn.

that study, we identified a 22-bp complemented palindrome from a highly conserved Coding Sequence (CDS). Both the 22-bp complemented palindrome (named Nankai complemented palindrome) and the CDS (named Nankai CDS), evolutionary conserved in BB coronavirus genomes, were identified as genomic features associated to the molecular functions of BB coronavirus. In the present study, we used these two genomic features to trace the origin of 2019-nCoV (GenBank: MN908947) and conduct a preliminary study of the mechanisms in the cross-species infection and host adaption of BB coronavirus. Our analytical results show that 2019-nCoV with large differences from the SARS coronavirus, may originate from BB coronaviruses in bats. The most important finding is that the alternative translation of Nankai CDS could produce more than 17 putative proteins, which may be responsible for the host adaption. The genotyping of 13 viruses using the 17 putative proteins revealed the high mutation rate and diversity of BB coronavirus. Our study, for the first time, aimed to explain the reason for the high host adaptability of the multi-host BB coronavirus at the molecular level using large amounts of genomic data. The findings in the present study laid foundations for the rapid detection, genotyping, vaccine development and drug design of, but not limited to BB coronavirus.

Keywords: Coronavirus; SARS; Alternative translation; SARS-CoV-2; Cross-species infection

冠状病毒属(*Coronavirus*)的病毒是具有包膜的正链 RNA 病毒。由于 2003 年严重急性呼吸综合征 (Severe Acute Respiratory Syndrome, SARS) 和 2012 年中东呼吸综合征 (Middle East Respiratory Syndrome, MERS) 的爆发,冠状病毒逐渐成为病毒学领域的一个研究热点。2019 年 12 月,中国武汉报道了冠状病毒引起的肺炎,其临床症状与 2003 年爆发的 SARS 不同,因此推断该病毒可能是冠状病毒的一个新变种。当前,冠状病毒溯源相对比较清晰:2003 年,Enserink 等在国际上首次提出了 SARS 冠状病毒 (SARS Coronavirus, SARS-CoV) 来自蝙蝠,而后通过果子狸传播给人类^[1];再后来的病毒溯源研究都是根据蝙蝠-果子狸-人途径回溯,并重点研究各种蝙蝠体内的 SARS 样 (SARS-like) 冠状病毒与已知 SARS 冠状病毒的进化关系;特别是 2013 年,SARS 冠状病毒 (GenBank: AY274119.3) 被溯源到中华菊头蝠 (*Rhinolophus sinicus*)^[2]。尽管已有足够的冠状病毒基因组数据公开,有关冠状病毒跨物种传播和宿主适应性方面的研究却一直难以突破,特别是无法解释冠状病毒变异快、宿主多且具有较强的宿主适应性等特点。其中一个最主要原因就是冠状病毒基因组较大,多个基因功能未知。例如,SARS 冠状病毒基因组 (GenBank: AY274119.3) 的全长是 29 751 bp,比 HIV-1 病毒全长三倍还要多。另外,冠状病毒基因组注释中基因命名不统一、且大部分注释的蛋白质序列都来自预测,无实验验证,难以保证正确性。高通量测序的迅速发展,使大量冠状病毒基因组序列得以公开。然而,当前大部分的冠状病毒基因组研究都是简单使用全基因组或某个病毒结构基因的序列,后一种做法主观性较大。简单使用全基因组序列进行进化分析和重要功能或致

病位点分析,不仅工作量巨大,而且大量与研究目的无关的基因组区域会对数据分析产生干扰。因此,不同于简单使用全基因组序列的其它研究,我们于 2018 年在国际上首次提出分子功能与进化分析相结合的研究思想,并应用于冠状病毒基因组的研究。分子功能与进化分析相结合,即仅使用与特定功能相关的基因组特征对一定变异范围内的序列进行进化分析^[3]。这样不仅可以简化进化分析、提高信噪比,使分析结果更可靠,而且可以用于指导功能实验。2018 年,南开大学高山等在国际上首次报道了 SARS 冠状病毒中存在互补回文小 RNA (complemented palindromic Small RNAs, cpsRNAs) 的现象,并确定了 Beta 冠状病毒 B 亚群 (BB 冠状病毒) 与功能相关的两个重要的基因组特征 (Genomic feature)^[3],即一段 22 bp 的互补回文序列 (命名为 Nankai complemented palindrome) 与其所在的编码区 (命名为 Nankai CDS)。在上述研究中,对 Nankai complemented palindrome 进行突变分析的结果与使用 Nankai CDS 进行进化分析的结果一致支持 SARS 冠状病毒的跨物种传播途径是蝙蝠-果子狸-人。在分子功能与进化分析相结合的研究思想的指导下,本研究使用 Nankai complemented palindrome 和 Nankai CDS 对新发布的 2019 新型冠状病毒基因组 (GenBank: MN908947) 进行分析以期准确溯源,并对 BB 冠状病毒跨物种传播和宿主适应性进行初步研究。

1 数据与方法

在前期研究中,我们通过互补回文小 RNA 的发现从 SARS 冠状病毒 MA-15 株 (GenBank:

DQ497008)中确定了BB冠状病毒中与功能相关的两个重要的基因组特征,即一段22 bp的互补回文序列与其所在的编码区。这个编码区(DQ497008:25689-26153)是BB冠状病毒基因组中一段高度保守的序列,在已知SARS冠状病毒基因组(如DQ497008)中只对应一个读码框ORF 3b;但在其它BB冠状病毒基因组中普遍对应多个读码框,因此ORF 3b无法代表其它BB冠状病毒基因组中的这段序列。本研究将22 bp的互补回文序列与其所在的编码区(可能是一个或多个读码框)分别命名为Nankai complemented palindrome和Nankai CDS。DQ497008和2003年爆发的SARS冠状病毒基因组(GenBank: AY274119.3)含有相同的Nankai complemented palindrome和Nankai CDS,因此在本研究中唯一代表SARS冠状病毒基因组。在前期研究中^[3],我们的研究对象仅限于BB冠状病毒(本研究亦是如此),并进行如下操作^[3]:提取NCBI GenBank数据库中全部BB冠状病毒的完整基因组(Complete genome)序列(以下简称病毒序列或序列);去除与Nankai CDS相似度在0.999以上的序列后,留下11条病毒序列(GenBank: JX993987、JX993988、GQ153539、GQ153540、GQ153542、DQ071615、DQ412042、DQ412043、AY515512、AY572034和DQ497008)用于进一步研究。本研究增加一条新发布的2019新型冠状病毒基因组的序列(GenBank: MN908947)和一条来自蝙蝠的序列(GenBank: MG772934),其Nankai CDS与2019新型冠状病毒Nankai CDS相似度最高(92.52%)。另外,由于本论文接收时2019新型冠状病毒基因组序列(GenBank: MN908947)尚未公开,实际使用的是来自互联网公开的序列(见致谢)。在本研究中,13条序列根据其宿主分为五组用于进一步研究,这五组命名为SARS(DQ497008)、果子狸Civet(AY515512和AY572034)、2019新型冠状病毒2019-nCoV(MN908947)、蝙蝠群体bat2(MG772934)和蝙蝠群体bat1(MG772934之外8条来自蝙蝠的序列)。在本研究中,进化分析使用软件PHYMLIP v3.69,统计与作图使用软件R v2.15.3^[4]。其它数据处理包括:病毒蛋白质预测和核苷酸三连子分析只考虑一个起始密码子ATG、三个终止密码子(TAA、TAG和TGA)和正向三个读码框;蛋白质预测不考虑小于20个氨基酸的蛋白质;最后,考虑四个起始密码子(ATG、GTG、CTG和TTG)和三个终止密码子重复蛋白质预测和核苷酸

三连子分析以检验结果的可靠性(见结果)。

2 结果

在经典定义中,DNA回文(DNA palindrome)和DNA互补回文(DNA complemented palindrome)序列都叫做DNA回文序列,而且大多数情况下,DNA回文序列的经典定义与我们新定义中的互补回文序列相同(例如ATCGCGAT)。根据我们的新定义^[3],回文序列(DNA或RNA)要求序列(一条链)从5'到3'方向与从3'到5'方向的读取结果相同(例如ATCGGCTA)。回文和互补回文序列应该加以区别,主要基于三个原因:(1)我们的前期研究发现了回文和互补回文序列可能具有不同的生物学意义;(2)在动物病毒基因组中,回文和互补回文序列具有不同的分布特征;(3)经典定义是从DNA水平定义的,我们的定义考虑了单链DNA和RNA(特别针对病毒)情况。2004年,Chew等在国际上首次报道了SARS冠状病毒基因组中较短互补回文序列(新定义)的分布特征^[5],其主要发现包括两点:(1)在所有分析的冠状病毒基因组中,长度为4 bp的互补回文序列出现频率显著较低;(2)长度为6 bp的互补回文序列仅在SARS冠状病毒(而不是所有冠状病毒)基因组中出现频率显著较低。因此,该研究的结论是SARS冠状病毒基因组含有较少的6 bp互补回文序列有可能利于该病毒有效规避宿主细胞内的某些防御机制而有利于病毒存活;该研究同时还发现了两个较长(14 bp以上)的互补回文序列分别是TCTTTAACAAGCTTGTAAAGA和TAAAATTAATTTTA。但是,仅仅根据概率模型计算得到的此类回文或互补回文序列数量巨大,绝大部分无法与分子功能关联,因此没有进一步研究的价值。直到2018年,我们偶然发现了互补回文小RNA恰好来自这个22 bp的互补回文序列(命名为Nankai complemented palindrome),特别是发现了根据这个序列预测的RNA二级结构在BB冠状病毒基因组中高度保守(见图1A)。RNA水平的证据是将这个序列与分子功能建立关联的关键^[3]。在前期研究中,我们对11个BB冠状病毒(见数据与方法)的Nankai complemented palindrome进行突变分析的结果与使用其所在的编码区(命名为Nankai CDS)进行进化分析的结果一致支持SARS冠状病毒的跨物种传播途径是蝙蝠-果子狸-人(见图1B红色方框内)^[3]。

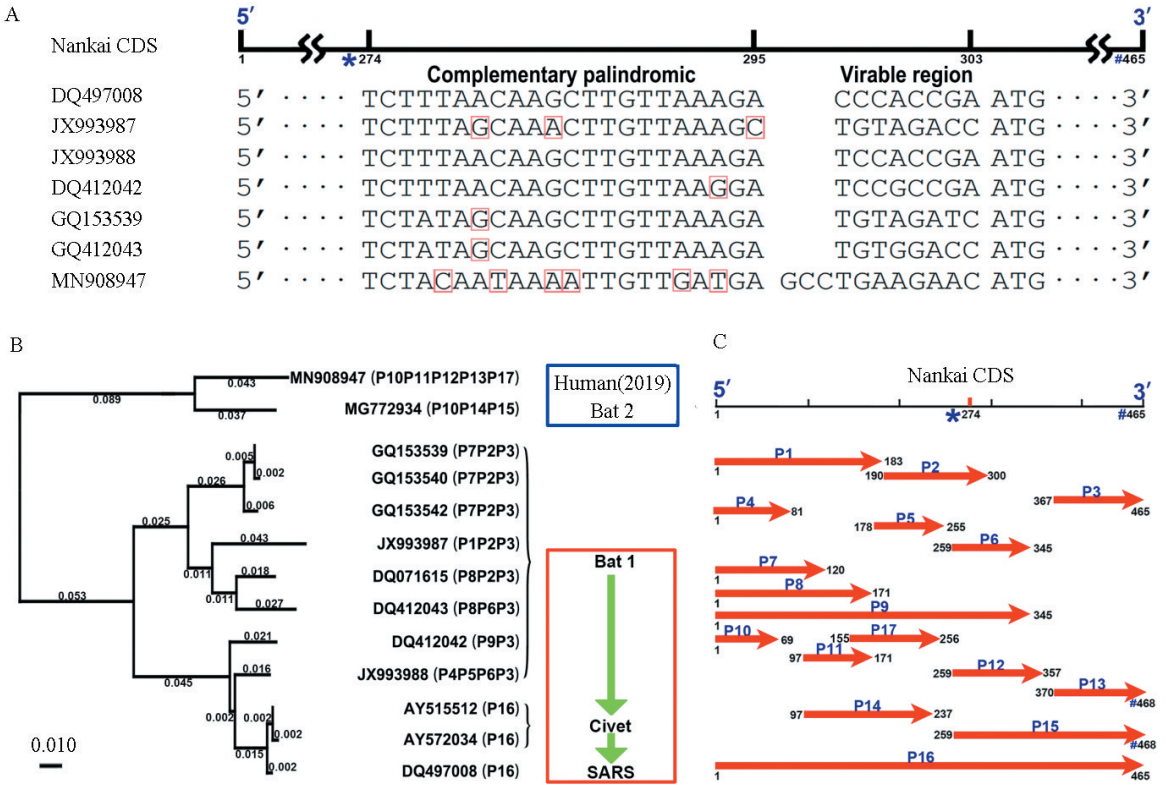


图 1 2019 新型冠状病毒的起源与可变翻译

Fig.1 Origin and alternative translation of the 2019 novel coronavirus genome

注:A:Nankai CDS 是 BB 冠状病毒基因组中一段高度保守的序列,包括一段 22 bp 的互补回文序列(用 * 表示)。Nankai CDS 还包括一个 8 或 11 bp 的可变区,分别属于两种不同长度(是 465 或 468 bp,用#表示)的 Nankai CDS;B :进化树构建使用 13 条去除可变区的 Nankai CDS。使用多种方法进行进化分析的结果一致,这里只展示邻接(Neighbor Joining, NJ)法的结果;2019 新型冠状病毒源自蝙蝠某个群体(蓝色方框内),但与 SARS 冠状病毒差异巨大(红色方框内);用 P1 到 P17(见图 1C)可以对 13 个 BB 冠状病毒进行基因分型(小括号内),其结果与进化分析的结果一致。C :可变翻译导致从 Nankai CDS 可以预测出至少 17 个蛋白质,分别命名为 P1 到 P17。这里展示的是使用一个起始密码子 ATG 和三个终止密码子进行蛋白质预测的结果,使用四个起始密码子(ATG、GTG、CTG 和 TTG)和三个终止密码子会预测更多的蛋白质,但不改变可变翻译相关结论。

Notes:A: The 22-bp complemented palindrome (named Nankai complemented palindrome) and the CDS (named Nankai CDS) have a high degree of evolutionary conservation in genomes of betacoronavirus subgroup B (BB coronavirus). An 8- or 11-bp long variable region in Nankai CDS resides in two types of Nankai CDS. The first type has a length of 465 bp, while the second type has a length of 468 bp. B : Consistent results were obtained by different phylogenetic analysis methods. Here, we only show the result using the NJ (Neighbor Joining) method. The 2019 novel coronavirus (in blue box) with large differences from the SARS coronavirus (in red box), may originate from BB coronaviruses in Chinese horseshoe bats. Phylogenetic analysis and genotyping (in parentheses) of 13 viruses using 17 putative proteins (Fig.1C) achieved consistent results. C: The putative proteins (named as P1 to P17) were predicted from the Nankai CDS using ATG as the start codon, and TAA, TAG and TGA as stop codons. The prediction using ATG, GTG, CTG and TTG as start codons, and TAA, TAG and TGA as stop codons did not change all conclusions in the present study.

在本研究中,所有 13 条 Nankai CDS(见数据与方法)都含有一个完整的 Nankai complemented palindrome 和一个紧邻其 3' 端的长度为 8 或 11 bp 的可变区 (Variable region)。使用 MN908947 (2019 新型冠状病毒)中 Nankai complemented palindrome 预测的 RNA 二级结构仍然能形成稳定的茎环 (Stem-loop) 结构,再次验证了 Nankai complemented palindrome 在 BB 冠状病毒基因组中的保守性^[3]。Nankai CDS 中的可变区导致了长度变异(见图 1A): MG772934 (来自蝙蝠)和 MN908947 (2019 新型冠状病毒)的 Nankai CDS 全长为 468 bp (含 11 bp 可变

区);其余 11 条序列的 Nankai CDS 全长为 465 bp (含 8 bp 可变区)。由于无法确定每条 Nankai CDS 的可变区的变异方式,进行进化分析前,13 条 Nankai CDS 去除可变区后即可对齐,长度统一到 465 bp。使用多种方法(UPGMA、邻接法、最大简约法、最大似然法和贝叶斯法)对 13 条 Nankai CDS 进行进化分析的结果(见图 1B)一致支持:(1) MN908947 (2019 新型冠状病毒)与 MG772934 (来自蝙蝠)在进化树的一个分支内(见图 1B 蓝色方框内),说明 2019 新型冠状病毒源自蝙蝠;(2) 2019 新型冠状病毒(见图 1B 蓝色方框内)与 SARS 冠状病毒(见图 1B 红色方

框内)在不同分支,而且距离很大,这一结果与两者临床症状差异一致;(3)SARS 冠状病毒的跨物种传播途径“蝙蝠-果子狸-人”(见图 1B 红色方框内)与 2018 年分析结果完全相同^[3],证明了我们方法的可靠性。

Nankai CDS 是 BB 冠状病毒基因组中一段高度保守的序列(见数据与方法),靠近其 3'端的 99 个碱基(编码 33 个氨基酸)极端保守,不仅没有起始或终止密码子突变,也极少有其它突变,具有重要的研究和应用(如引物设计)价值。Nankai CDS 在来自 SARS 和果子狸的 beta 冠状病毒基因组中是一个完整的读码框,编码蛋白 P16(见图 1C),但在来自其它宿主的冠状病毒基因组中(不包括 3'端的 99 个碱基)出现了大量的起始或终止密码子突变(见表 1),并导致 Nankai CDS 出现读码框的多样性,即可变翻译。由于去冗余后的完整基因组数据较少,需要从两个方面排除这些突变是来自测序错误等导致的假阳性。一方面,用 13 条 Nankai CDS 比对到 NCBI NT 数据库,看是否有其它冠状病毒序列支持这些突变。比对结果显示,大部分序列上的突变都能得到 NT 数据库中其它数据支持,只有三个例外:JX993988 在

Nankai CDS 上的三个突变(C80A、T80A 和 T254A),MG772934(来自蝙蝠)和 MN908947(2019 新型冠状病毒)在 Nankai CDS 上共有的四个突变(T182A、C182A、T237A 和 T236G)没有其它数据支持。JX993988 代表的病毒类型比较特殊,与其它来自蝙蝠的病毒差异较大(见图 1B),这里不做进一步研究。另外一方面,统计 13 条 Nankai CDS 中核苷酸三连子出现次数,次数最高的前 14 位依次是:ATG(123 次)、ACT(106 次)、AAG(98 次)、TTG(83 次)、CAG(77 次)、CAA(71 次)、GTG(68 次)、AAA(66 次)、CTA(65 次)、ACA(60 次)、TTA(59 次)、ATT(58 次)、CGA(57 次)和 CTG(55 次)。核苷酸三连子统计结果显示:前 14 位三连子出现次数之和超过三连子总数一半以上(1 046 次/2 016 次);6 个三连子与起始密码子 ATG 只有 0 或 1 个核苷酸的差异;8 个三连子与终止密码子(TAA、TAG 和 TGA)只有 1 个核苷酸的差异。统计结果说明了通过 1 个核苷酸突变(点突变)导致起始或终止密码子的丢失和获得的概率非常高。根据以上两个方面的分析结果,基本可以排除 Nankai CDS 中的大量起始或终止密码子突变是假阳性(见表 1)。

表 1 Nankai CDS 内起始或终止密码子突变

Table 1 Mutations in start or stop codons in Nankai CDS

PP	Start	Allele	Mutation	Stop	Allele	Mutation	AA
P1	1	ATG		183	TTG\TCG\TAG	T182A\C182A	61
P2	190	ATA\ATG	A192G	300	TGG\TAG	G299A	37
P3	367	ATG		465	TAA		33
P4	1	ATG		81	TCG\TTG\TAG	C80A\T80A	27
P5	178	GTG\ATG	G178A	255	TTG\TAG	T254A	26
P6	259	ACG\ATG	C260T	345	TTA\TAA	T344A	29
P7	1	ATG		120	TCA\TTA\TAA	C119A\T119A	40
P8	1	ATG		171	TTA\TAA	T170A	57
P9	1	ATG		345	TTA\TAA	T344A	115
P10	1	ATG		69	TCA\TAA	C68A	23
P11	97	ACG\ATG	C98T	171	TTA\TAA	T170A	25
P12	259	ACG\ATG	C260T	357	CAA\TAA	C355T	33
P13	370	ATG	468	TAA	33		
P14	97	ACG\ATG	C98T	237	TGT\TTA\TGA	T237A\T236G	47
P15	259	ACG\ATG	C260T	468	TAA	70	
P16	1	ATG	465	TAA	155		
P17	155	GTG\TTG\ATG	G155A\T155A	256	TCA\AGA\TGA	C255G\A254T	34

注:可变翻译导致从 Nankai CDS 可以预测出至少 17 个蛋白质,分别命名为 P1 到 P17(见图 1C)。所有突变均来自 17 个预测的蛋白质,PP 表示预测的蛋白质(Putative Protein)的名称;Start 一列记录的是起始密码子的第一个碱基的位置;Stop 一列记录的是终止密码子的最后一个碱基位置;Allele 一列记录所在位置的密码子的全部类型;Mutation 一列只记录导致获得起始或终止密码子的突变,突变格式为“突变前碱基-突变位置-突变后碱基”;AA 表示 PP 的长度,单位是氨基酸(Amino Acid)。这里展示的是只使用一个起始密码子 ATG 和三个终止密码子进行蛋白质预测的结果;使用四个起始密码子(ATG、GTG、CTG 和 TTG)和三个终止密码子会预测更多的蛋白质,但不改变可变翻译相关结论。

Notes: The putative proteins (named as P1 to P17) were predicted from the Nankai CDS (Fig.1C). This table only includes the mutations in start or stop codons in 17 putative proteins. "PP" is Putative Protein; "Start" records the positions of the first nucleotides in start codons; "Stop" records the positions of the last nucleotides in stop codons; "Allele" records all types of codons on the positions; "Mutation" only records mutations which cause the acquisition of start or stop codons in the format "nucleotide before mutation-position-nucleotide after mutation"; "AA" is Amino Acid. The putative proteins were predicted from the Nankai CDS using ATG as the start codon, and TAA, TAG and TGA as stop codons. The prediction using ATG, GTG, CTG and TTG as start codons, and TAA, TAG and TGA as stop codons did not change all conclusions in the present study.

可变翻译导致从 Nankai CDS 可以预测出至少 17 个蛋白质 (Putative Protein), 分别命名为 P1 到 P17 (见图 1C)。用 P1 到 P17 可以对 13 个病毒进行基因分型 (基于可变翻译的基因分型), 例如, MN908947 (2019 新型冠状病毒) 可以表示为 P10P11P12P13P17。我们发现基因分型相同或相近的病毒都在进化树上的同一个分支内 (见图 1B)。例如, bat1 分支的病毒的基因分型都由三个蛋白质组成; SARS 和果子狸分支都是 P16; MN908947 (2019 新型冠状病毒) 和 MG772934 (来自蝙蝠) 都有 P10。根据以上事实, 可以得到以下推论: (1) BB 冠状病毒 (不限于 Nankai CDS 内) 存在大量的可变翻译; (2) BB 冠状病毒变异快、多样性高, 仅仅 13 个病毒即可得到多达 9 种基因分型 (见图 1BC); (3) 在 SARS 冠状病毒从蝙蝠到人传播过程中, Nankai CDS 内多个点突变积累才形成 P16; (4) 在 2019 新型冠状病毒从蝙蝠到人传播过程中, P15 只需一个点突变 (C355T) 即形成 P12 和 P13, P14 只需两个点突变 (T170A 和 T191C) 即形成 P11 和 P17; (5) 可变翻译可用于 (但不限于) BB 冠状病毒的基因分型, 用 P1 到 P17 进行基因分型的结果与用 Nankai CDS 进行进化分析的结果一致。 (6) BB 冠状病毒可能通过可变翻译以适应不同宿主。

3 结 论

本研究证明了我们提出的分子功能与进化分析相结合的研究思想的可行性与可靠性, 对生物大分子的功能或进化研究具有一定的指导意义, 并对病毒或细菌等病原体的研究提供了新的路线和方向。因此, 本研究除了解决 2019 新型冠状病毒溯源等具体问题, 更主要贡献是提供了研究思想和方法。

1) 使用 Nankai CDS 溯源分析的结果支持 2019 新型冠状病毒源自蝙蝠, 但与 SARS 冠状病毒差异巨大;

2) 从 BB 冠状病毒可变翻译中获取的信息可应用于其快速检测、基因分型、疫苗开发以及药物设计;

3) BB 冠状病毒可能通过可变翻译以适应不同宿主;

4) 对 BB 冠状病毒可变翻译的研究将有助于研

究该病毒的变异、宿主适应性、感染以及致病机制等问题。

4 讨 论

病毒或细菌等原核生物的可变翻译现象很早就已经发现, 但是其生物学意义尚未清楚, 特别是缺乏基于较大数据的实证研究。我们在前期研究植物病毒^[6]、昆虫病毒^[7]、人病毒^[8]、哺乳动物病毒 (特别是非洲猪瘟病毒^[9]) 和布鲁氏菌^[10] 的过程中, 都没有明显观察到可变翻译现象, 特别是 BB 冠状病毒 Nankai CDS 中这样大量而集中存在的可变翻译 (仅仅 13 个病毒即可得到多达 9 种基因分型) 极其罕见。此次实证研究结果达到了基于可变翻译的基因分型与进化分析的一致性, 是一次较大突破, 主要得益于以下几点: (1) 分子功能与进化分析相结合的研究思想得以应用; (2) 有蝙蝠-果子狸-人这一跨物种传播途径作为方向性指导; (3) 无意中发现 Nankai complemented palindrome 和 Nankai CDS; (4) 来自多种宿主 (蝙蝠、果子狸和人) 的病毒数据, 特别是阳性样本 (如来自 SARS 和武汉肺炎患者) 数据的积累。经典理论认为病毒与宿主共用一套翻译系统, 然而病毒的蛋白质翻译也有其特殊性, 导致本研究从 Nankai CDS 中预测的 17 种蛋白质不能提供准确的序列。蛋白质序列的准确度不影响本研究的所有结论和推论, 而且为下一步蛋白质、细胞和动物水平的实验指出方向。

基于可变翻译的基因分型相同的病毒都在进化树上的同一个分支 (同种宿主) 内预示了 beta 冠状病毒可能通过可变翻译以适应不同宿主 (见图 1B)。但是, 由于目前能获得的病毒数据仅仅来自少量宿主 (蝙蝠、果子狸和人), 这一推断的验证还需要积累来自更多种类宿主的病毒的数据。我们推测, 病毒通过可变翻译以适应不同宿主作为一种普遍机制, 并不限于冠状病毒或原核生物; 真核生物中某些基因的可变剪接可能也有相似或相近的生物学功能。

BB 冠状病毒通过可变翻译产生的多种蛋白质的致病性研究是一个更为重要的课题。当前公开的基因组数据虽然总量很大, 但是, 实际有效的阳性样本 (有实验或临床记录证明其致病性) 种类不足和缺

乏阴性(临床实验确认的非 SARS 病毒)样本。致病性最明确的 2019 新型冠状病毒和 SARS 冠状病毒在 Nankai CDS 上各产生五个小蛋白和一个大蛋白(见图 1C),并没有显示出蛋白质某个属性(如大小)与致病性的关系。在预测的 17 个蛋白质中,只有 P16(长度为 155 个氨基酸)和 P9(长度为 115 个氨基酸)是大蛋白质(见表 1);P16(来自 SARS 冠状病毒)的致病性已知,而 P9(来自蝙蝠的冠状病毒)未知。因此,来自蝙蝠的其它冠状病毒产生的蛋白(如 P9)未必没有致病性,或许也具有导致 2019 肺炎或 SARS 的能力。本研究预测的 17 种蛋白质大小适中而且差异明显,可以很容易通过蛋白质组或外源表达等技术进一步验证其致病性。

致谢:感谢新型冠状病毒相关工作人员(省疾控中心等)的基础性工作,特别复旦大学张永振等公开的 2019 新型冠状病毒基因组数据,感谢南开大学生命科学学院动物系卜文俊、张涛、黄大卫、刘燕强和贺秉军等各位老师对我们生物信息学研究的长期支持。

参考文献(References)

- [1] ENSERINK M, NORMILE D, VOGEL G. Clues to the animal origins of SARS[J]. *Science*, 2003, 300(5624): 1351–1351. DOI: 10.1126/science.300.5624.1351a.
- [2] GE Xingyi, LI Jialu, YANG Xinglou, et al. Isolation and characterization of a bat SARS-like coronavirus that uses the ACE2 receptor[J]. *Nature*, 2013, 503(7477): 535–538. DOI: 10.1038/nature12711.
- [3] LIU Chang, CHEN Ze, HU Yue, et al. Complemented Palindromic Small RNAs First Discovered from SARS Coronavirus[J]. *Genes*, 2018, 9(9): 1–11. DOI: 10.3390/genes9090442.
- [4] GAO Shan, OU Jianhong, XIAO Kai. *R language and Bioconductor in bioinformatics applications (Chinese Edition)* [M]. Tianjin: Tianjin Science and Technology Translation Publishing Ltd. 2014.
- [5] CHEW D S H, CHOI K P, HEIDNER H, et al. Palindromes in SARS and other coronaviruses[J]. *Inform Journal on Computing*, 2004, 16(4): 331–340. DOI: 10.1287/ijoc.1040.0087.
- [6] LI Rugang, GAO Shan, HERNANDEZ A G, et al. Deep sequencing of small RNAs in tomato for virus and viroid identification and strain differentiation[J]. *PLoS One*, 2012, 7(5): 1–10. DOI: 10.1371/journal.pone.0037127.
- [7] CHEN Yunru, ZHONG Silin, FEI Zhangjun, et al. Transcriptome responses of the host *trichoplusia ni* to infection by the baculovirus *autographa californica* multiple nucleopolyhedrovirus[J]. *Journal of Virology*, 2014, 88(23): 13781–13797. DOI: 10.1128/JVI.02243–14.
- [8] WANG Fang, SUN Yu, RUAN Jishou, et al. Using small RNA deep sequencing data to detect human viruses[J]. *BioMed Research International*, 2016, 2016(2016): 1–9. DOI: 10.1155/2016/2596782.
- [9] CHEN Ze, XU Xiaofeng, WANG Yufeng, et al. DNA segments of African Swine Fever Virus detected for the first time in hard ticks from sheep and bovines[J]. *Systematic & Applied Acarology*, 2019, 24(1): 180–184. DOI: 10.11158/saa.24.1.13.
- [10] WANG Yuanzhi, WANG Zhen, CHEN Xin, et al. The complete genome of *brucella suis* 019 provides insights on cross-species infection[J]. *Genes*, 2016, 7(2): 1–12. DOI: 10.3390/genes7020007.

[责任编辑:吴永英]