

LA-UR- 09-01624

Approved for public release;  
distribution is unlimited.

Title: Genome Assortment, Not Serogroup,  
Defines *Vibrio cholerae* Pandemic Strains

Author(s): Jongsik Chun, Christopher J. Grim, Nur A. Hasan, Je Hee Lee, Seon Young Choi, Bradd J. Haley, Elisa Taviani, Yoon-Seong Jeon, Dong Wook Kim, Jae-Hak Lee, Thomas S. Brettin, David C. Bruce, Jean F. Challacombe, J. Chris Detter, Cliff S. Han, A. Christine Munk, Olga Chertkov, Linda Meincke, Elizabeth Saunders, Ronald A. Walters, Anwar Huq, G. Balakrish Nair and Rita R. Colwell

Intended for: Nature



Los Alamos National Laboratory, an affirmative action/equal opportunity employer, is operated by the Los Alamos National Security, LLC for the National Nuclear Security Administration of the U.S. Department of Energy under contract DE-AC52-06NA25396. By acceptance of this article, the publisher recognizes that the U.S. Government retains a nonexclusive, royalty-free license to publish or reproduce the published form of this contribution, or to allow others to do so, for U.S. Government purposes. Los Alamos National Laboratory requests that the publisher identify this article as work performed under the auspices of the U.S. Department of Energy. Los Alamos National Laboratory strongly supports academic freedom and a researcher's right to publish; as an institution, however, the Laboratory does not endorse the viewpoint of a publication or guarantee its technical correctness.

# **Genome Assortment, Not Serogroup, Defines *Vibrio cholerae* Pandemic Strains**

Jongsik Chun <sup>1,2,3</sup>, Christopher J. Grim <sup>2</sup>, Nur A. Hasan <sup>4,6</sup>, Je Hee Lee <sup>1,3</sup>, Seon Young Choi <sup>1,3</sup>, Bradd J. Haley <sup>4</sup>, Elisa Taviani <sup>4</sup>, Yoon-Seong Jeon <sup>3</sup>, Dong Wook Kim <sup>3</sup>, Jae-Hak Lee <sup>1</sup>, Thomas S. Brettin <sup>7</sup>, David C. Bruce <sup>7</sup>, Jean F. Challacombe <sup>7</sup>, J. Chris Detter <sup>7</sup>, Cliff S. Han <sup>7</sup>, A. Christine Munk <sup>7</sup>, Olga Chertkov <sup>7</sup>, Linda Meincke <sup>7</sup>, Elizabeth Saunders <sup>7</sup>, Ronald A. Walters <sup>8</sup>, Anwar Huq <sup>4</sup>, G. Balakrish Nair <sup>9</sup> and Rita R. Colwell <sup>2,4,5</sup>

<sup>1</sup>, School of Biological Sciences and Institute of Microbiology, Seoul National University, Seoul 151-742, Republic of Korea,

<sup>2</sup>, Center of Bioinformatics and Computational Biology, University of Maryland Institute of Advanced Computer Studies, University of Maryland, College Park, MD 20742, U.S.A.

<sup>3</sup>, International Vaccine Institute, Seoul, 151-818, Republic of Korea,

<sup>4</sup>, Maryland Pathogen Research Institute, University of Maryland, College Park, MD 20742, USA

<sup>5</sup>, Johns Hopkins Bloomberg School of Public Health, Baltimore, Maryland, USA

<sup>6</sup>, International Center for Diarrheal Disease Research, Bangladesh, Dhaka-1000, Bangladesh

<sup>7</sup>DOE Joint Genome Institute, Bioscience Division, Los Alamos National Laboratory, Los Alamos, NM 87545, USA

<sup>8</sup>, Pacific Northwest National Laboratory, Richland, WA 99352, USA

<sup>9</sup>, National Institute of Cholera and Enteric Diseases, Beliaghata, Kolkata 700 010, India

Correspondence and requests for materials should be addressed to Rita R. Colwell (rcolwell@umiacs.umd.edu)

*Vibrio cholerae*, the causative agent of cholera, is a bacterium autochthonous to the aquatic environment, and a serious public health threat. *V. cholerae* serogroup O1 is responsible for the previous two cholera pandemics, in which classical and El Tor biotypes were dominant in the 6th and the current 7th pandemics, respectively. Cholera researchers continually face newly emerging and re-emerging pathogenic clones carrying combinations of new serogroups as well as of phenotypic and genotypic properties. These genotype and phenotype changes have hampered control of the disease. Here we compare the complete genome sequences of 23 strains of *V. cholerae* isolated from a variety of sources and geographical locations over the past 98 years in an effort to elucidate the evolutionary mechanisms governing genetic diversity and genesis of new pathogenic clones. The genome-based phylogeny revealed 12 distinct *V. cholerae* phyletic lineages, of which one, designated the *V. cholerae* core genome (CG), comprises both O1 classical and El Tor biotypes. All 7th pandemic clones share nearly identical gene content, i.e., the same genome backbone. The transition from 6th to 7th pandemic strains is defined here as a "shift" between pathogenic clones belonging to the same O1 serogroup, but from significantly different phyletic lineages within the CG clade. In contrast, transition among clones during the present 7th pandemic period can be characterized as a "drift" between clones, differentiated mainly by varying composition of laterally transferred genomic islands, resulting in emergence of variants,

**exemplified by *V.cholerae* serogroup O139 and *V.cholerae* O1 El Tor hybrid clones that produce cholera toxin of classical biotype. Based on the comprehensive comparative genomics presented in this study it is concluded that *V. cholerae* undergoes extensive genetic recombination *via* lateral gene transfer, and, therefore, genome assortment, not serogroup, should be used to define pathogenic *V. cholerae* clones.**

*Vibrio cholerae*, a bacterium autochthonous to the aquatic environment, is the causative agent of cholera, a severe, watery, life-threatening diarrheal disease. Historically, cholera bacteria have been serogrouped based on their somatic O antigens, with more than 200 serogroups identified to date <sup>1</sup>. Although strains from many of the serogroups of *V. cholerae* have caused either individual cases of mild gastroenteritis or local outbreaks of gastroenteritis, only the toxigenic strains of serogroups O1 and O139 have been identified as agents of cholera epidemics. Genes coding for cholera toxin, *ctxAB*, and other virulence factors have been shown to reside in bacteriophages and various mobile genetic elements. In addition, *V. cholerae* serogroup O1 is differentiated into two biotypes, classical and El Tor, by a combination of biochemical traits and sensitivity to specific bacteriophages <sup>2</sup>.

Throughout human history cholera pandemics have been recorded with seven such pandemics characterized over the past hundred or more years. Today the disease remains endemic only in developing countries, even though *V. cholerae* is native to estuaries and river systems throughout the world <sup>3</sup>. Isolates of the 6th pandemic were almost exclusively of the O1 classical biotype, whereas the current (7th) pandemic is dominated by *V. cholerae* O1 El Tor biotype as the causative agent, a transition occurring between 1905 and 1961. The six pandemics previous to the current pandemic are considered to have originated in the Indian subcontinent, whereas the 7th pandemic strain was first isolated in the Indonesian island of

Sulawesi in 1961, and subsequently in Asia, Africa, and Latin America.

Over the last twenty years, several new epidemic lineages of *V. cholerae* O1 El Tor have emerged or re-emerged. In 1992, a new serogroup of *V. cholerae*, O139, was identified as the cause of epidemic cholera in India and Bangladesh <sup>4</sup>. That is, both *V. cholerae* O1 El Tor and O139 consistently have been isolated where the major cholera epidemics since 1992 have occurred, although the geographic location of *V. cholerae* O139 appears still to be restricted to Asia. Additionally, *V. cholerae* "hybrid" O1 El Tor variants that carry the classical type CTX prophage, or produce classical type cholera toxin subunit B have been repeatedly isolated in Bangladesh <sup>5,6</sup> and Mozambique <sup>7</sup>. These new variants have subsequently replaced the prototype 7th pandemic *V. cholerae* O1 El Tor strains in Asia and Africa, with respect to frequency of isolation from clinical cases of cholera.

It is clear that the dynamics of *V. cholerae*, like other enteric pathogens, can be attributed to extensive lateral gene transfer *via* transduction, conjugation, and transformation <sup>2,8,9</sup>. However, the evolutionary history of this bacterium remains to be documented. Here we compare the genome sequences of 23 *V. cholerae* strains, representing diverse serogroups that have been isolated at various times over the past 98 years from a variety of sources and geographical locations. We conclude that the current pandemic is caused by strains belonging to a single phyletic line, diversified mainly by lateral gene transfer occurring in the natural environment.

### **Phylogeny and gene content of *Vibrio cholerae*.**

Phylogenetic analysis, accomplished using *ca.* 1.4 million bp of orthologous protein-coding regions for 23 *V. cholerae* strains (Table 1), revealed 12 distinct phyletic lineages. Strains belonging to non-O1/non-O139 serogroups from various sources showed substantial genomic diversity (Fig. 1a). In contrast, all *V. cholerae* serogroup O1 strains, except for two, comprised a monophyletic clade, designated *V. cholerae* core genome (CG) clade. Strains of both the 6th and 7th pandemics are concluded to have evolved from a common ancestor of this CG clade.

Twelve strains of the CG clade were further divided into two subgroups, as shown in the phylogenetic tree constructed using *ca.* 2.6 million bp alignment (Fig. 1a,b). The CG-1 subclade is comprised of most of the *V. cholerae* O1 El Tor strains and one *V. cholerae* O139 strain, whereas the CG-2 subclade contains strains of *V. cholerae* O1 classical and O37 serogroups. Interestingly, all clinical isolates associated with the current 7th cholera pandemic formed a very tight, monophyletic clade within the CG-1 subclade, which we have designated the 7th pandemic (7P) clade (Fig. 1b). *V. cholerae* O1 El Tor and O139 strains isolated from the Indian subcontinent and Africa epidemics during 1975 to 2004 are located in the 7P clade.

*V. cholerae* O1 El Tor N16961, the first *V. cholerae* strain to have been sequenced<sup>10</sup>,

contains 3,908 ORFs, according to the RAST annotation server <sup>11</sup>. Based on a comprehensive ortholog detection using reciprocal comparison, we identified 2,432 "core" ORFs (66.9%) in all 23 *V. cholerae* strains, results similar to those reported for *Escherichia coli* (2,344 <sup>12</sup> and 2,865 <sup>13</sup>). Core gene ratios of the large and small chromosomes were 69.7 and 59.9% of the total genome of strain N16961, respectively. The higher core gene content obtained for the large chromosome implies a functional importance for genes in this chromosome. Heidelberg et al. <sup>10</sup> had shown that the large chromosome contained fewer genes coding for hypothetical proteins than the small chromosome.

Comparative genomic analysis was used to identify a set of non-redundant orthologues, yielding 6,953 gene sets of the pan-genome of *V. cholerae* based on 23 strains. The number of non-redundant genes found in 7P, CG-1 and CG clades were 4,327, 4,664 and 4,818, respectively. Overall, more genes were newly found when two strains showed further evolutionary distance (Supplementary Fig. 1b). On average, an increase of 206 new genes occurred each time a strain was added to the *V. cholerae* gene pool from which the pan-genome was calculated, if the strains that were added represented unique phyletic lines outside the CG group (Supplementary Fig. 1c). The number of additional genes per strain added is significantly higher than reported for *Streptococcus agalactiae* (27 genes) <sup>14</sup>, but less than for *E. coli* (300 genes) <sup>12</sup>.



### **O serogrouping in the context of genome evolution.**

The lipopolysaccharide (LPS) of *V. cholerae* consists of three major regions: lipid A, core oligosaccharide (OS), and O antigen. *V. cholerae* synthesizes core OS and O antigen using *wav* and *wb\** gene clusters, respectively<sup>15,16</sup>. Molecular phylogeny and genetic organization of the *wav* and *wb\** gene clusters, located adjacently in the large chromosome, are summarized in Fig. 2. Nesper *et al.*<sup>15</sup> surveyed *wav* gene clusters of 38 *V. cholerae* strains and categorized them into five types and, in this study, core OS types 1, 3, 4 and 5 were detected. Interestingly, *V. cholerae* biovar albensis VL426 contains a variant of type 3, which we have designated type 3b.

In contrast to the limited diversity observed in the *wav* gene cluster (5 major types), 11 different types of *wb\** gene clusters were observed among the 23 strains. Phylogeny and genetic organization, based on the whole genome (Fig. 1), core OS, and O antigen gene clusters (Fig. 2a), clearly indicate both core OS and O antigen gene clusters are mobile *via* lateral gene transfer. The relatively stable gene order (synteny) of the core OS gene cluster indicates that it transfers as an entity. In contrast, the region coding for the O antigen is comprised of combinations of several smaller gene sets of different origin, leading to a remarkable diversity of the various O antigens seen in nature (Fig. 2b). This finding is in good agreement with a previous study<sup>17</sup> showing that the gene cluster coding for the O139 antigen is similar in *V. cholerae* serogroup O22, where substitution of a part of the cluster

occurred, but not a deletion.

Genome phylogeny (Fig. 1a) revealed that strains of O1 serogroup are found in three different phyletic lineages, namely the CG clade, and the *V. cholerae* O1 El Tor 12129(1) and TM11079-80 strains, in which the coding region for the O1 antigen is nearly identical. It is concluded that the O1 antigen phenotype arose at least three times in the evolution of *V. cholerae*. Furthermore, we hypothesize that the ancestor of the CG clade possessed a combination of the type 1 core OS and the O1 antigen gene clusters, giving rise to the present 12 *V. cholerae* CG strains, including the two *V. cholerae* non-O1 strains (V52 and MO10). The latter two became different serogroups by gene replacement, *via* lateral gene transfer, with strain V52 receiving both type 1 core OS and O37 antigen gene clusters from a *V. cholerae* O37 strain and *V. cholerae* MO10 receiving only the *V. cholerae* O139 antigen gene cluster from an unknown source, possibly a variant of the *V. cholerae* O22 serogroup<sup>17</sup>.

The *V. cholerae* O1 strains not belonging to the CG group, *V. cholerae* 12129(1) and TM11079-80, are environmental isolates from Australia isolated in 1985<sup>18</sup>, and from Brazil, isolated in 1980<sup>19</sup>. They showed the typical El Tor phenotype, but unlike other *V. cholerae* O1 El Tor strains in the CG-1 subclade, lack the two major virulence-related genomic islands, i.e., CTX prophage containing cholera toxin genes (*ctxAB*) and *Vibrio* pathogenicity island-1 (VPI-1) containing genes for biosynthesis of the toxin co-regulated pilin (TCP). By comparing genome phylogenies based on the whole genome (Fig. 1) and gene clusters coding

for the core OS (Fig. 2a) and O1 antigen (Fig. 2c), it is clear that genesis of these non-toxicogenic *V. cholerae* O1 El Tor strains can be attributed to independent lateral gene transfer events, perhaps transfer of only the O1 antigen gene cluster, but not the core OS region. The O1 antigen gene cluster itself appears to be highly conserved, since only 66 nucleotides were found to vary within the 19,503 bp multiple sequence alignment spanning 19 orthologues among the 11 *V. cholerae* O1 strains sequenced in this study. In any case, the phylogenetic analysis suggests, that donors of the O1 antigen gene cluster of *V. cholerae* 12121(1) and TM11079-80 were most likely members of the CG-2 subclade encompassing *V. cholerae* O1 classical biotype (Fig. 2c).

Thus, at least four O serogroup conversions, were detected, i.e., conversion from non-O1 to O1 (twice), O1 to O139, and O1 to O37, among the 23 *V. cholerae*. Earlier, several non-genomics based studies suggested such conversions take place in nature <sup>16,19,20</sup>, and chitin-induced natural transformation has been proposed as a mechanism in the natural environment <sup>21</sup>. Subsequently, *V. cholerae* O1 to O139 serogroup conversion by a single-step exchange of large fragments of DNA was demonstrated in a microcosm experiment <sup>9</sup>, clearly supported by the conclusion of this study that O serogroup conversion occurs frequently in nature. Mobility of the O phenotype in *V. cholerae* was first proposed by Colwell et al. <sup>22</sup>, and the cumulative results of both *in vivo* and *in vitro* experiments, as well as the genomic evidence, are compelling. Given the inconsistency between O serogroup typing and genome-

based phylogeny, it is concluded that, at the very least, the term "O1 El Tor" is both misleading and inaccurate for describing a set of phylogenetically coherent *V. cholerae* strains, especially in light of the frequency of serogroup conversion. Therefore, we propose a new terminology based on genome sequence; namely the core genome (CG) clade, CG-1 and CG-2 subclades, and 7th pandemic (7P) clade, to describe homologous intraspecific groups of *V. cholerae* (Fig. 1).

#### **Virulence-associated prophage and genomic islands within the context of the genome.**

*V. cholerae* possesses several known virulence factors, of which the cholera toxin (CT) and toxin-coregulated pilus (TCP) are considered to be the most significant. Genes coding for CT (*ctxAB*) are part of a temperate filamentous bacteriophage CTX $\phi$ <sup>8</sup> that can be incorporated into both chromosomes of *V. cholerae* at specific positions. The CTX $\phi$  genes were found to be present in members of the CG clade, except for *V. cholerae* NCTC 8457 and 2740-80. Among strains not belonging to the CG clade, only *V. cholerae* serogroup O141 (V51) contains this prophage.

The CTX $\phi$  found in classical and El Tor biotypes differs in the sequence of their repressor gene, *rstR*, and are classified as CTX $\phi$ <sup>Class</sup> and CTX $\phi$ <sup>El Tor</sup>, according to the biotype of the host<sup>23</sup>. From the genome sequences, it was found that CTX $\phi$ <sup>Class</sup> is not restricted to the classical biotype, but is also widely distributed in *V. cholerae* O1 El Tor and O141 strains

(Fig. 3). Given that *V. cholerae* O1 El Tor MAK 757, a clinical strain isolated in 1937, has this type of prophage, correlation of host biotype and prophage type is not considered significant. A recent study<sup>24</sup> showing the infection of CTX $\phi$ <sup>Class</sup> to *V. cholerae* non-O1 supports our finding.

Chromosomal attachment sites for CTX $\phi$  are known to harbor other genetic elements, including toxin-linked cryptic (TLC), RS1 elements, and VSK(=pre-CTX) prophages<sup>25,26</sup>. Additionally, we have discovered five genomic islands (GI-19, GI-22, GI-33, GI-43, GI-48) in the region of the CTX $\phi$  attachment sites on both chromosomes. In total, nine distinct genetic elements were found in these regions, where they appear in different combinations (Fig. 3). Seven strains possess GI-19 in either chromosome, which is similar but not identical to KSF-1 $\phi$ , discovered in an environmental *V. cholerae* strain<sup>27</sup>. It is evident that more bacteriophages/genetic elements are located in the CTX $\phi$  attachment regions of CG strains than non-CG strains. The ability to harbor more, especially the toxigenic, bacteriophage-like elements, in these regions in CG strains may explain why only members of the CG strains have been agents of cholera pandemics over the last century, namely during the 6th and 7th pandemics. We found no two toxigenic (CTX $\phi$ -harboring) strains with an identical GI organization and combination, except for two "hybrid" strains (the only 7P members harboring CTX $\phi$  |<sup>Class</sup>). It is evident from Fig. 3 that the two CTX $\phi$  attachment sites serve as an engine of genetic diversity for *V. cholerae* CG clade.

Genes coding for the toxin-coregulated pilus (TCP) are part of a genomic island, VPI-1, that is present in all CG strains. Among the non-CG strains, only *V. cholerae* O141 V51 contained VPI-1 but with less sequence similarity. Since TCP serves as a receptor for CTX $\phi$ , it explains why only this strain, among all of the non-CG strains, possesses CTX $\phi$ . Results of phylogenetic analysis based on the 24 genes of VPI-1 suggest that the original GI of *V. cholerae* NCTC 8457 has subsequently been replaced by a VPI-1 of a non-CG strain (Supplementary Fig. 2). Interestingly, GI-47, but not VPI-1, was found in the same genomic region in strains MZO-3, 1587, MZO-2, and VL426. This cassette-like property of GI mobility was also observed for the other known pathogenicity islands, including VPI-2, VSP-1, and VSP-2 (Supplementary Table 3).

### **Extensive lateral gene transfer in *V. cholerae*.**

Since it is generally accepted that lateral gene transfer plays an important role in the evolution of pathogenic bacteria, *V. cholerae* can serve as a useful paradigm. In this study, we define a GI as a genomic region containing five or more ORFs, where transfer, but not deletion, is obvious from comparison of the genome phylogeny and its presence/absence among test strains. A total of 73 GIs were identified (Supplementary Table 2) and their chromosomal locations are shown in Fig. 4. As discussed earlier, with respect to GIs associated with O antigen biosynthesis, CTX $\phi$ , VPI-1,2 and VSP-2, a total of 13 (eight in the

large and five in the small chromosome) genomic regions were found to have a cassette-like property, whereby different GIs occupy the same or similar region (Supplementary Table 3). Most GIs were found to be singletons in a given genome, although two (GI-12, GI-21) were present as four and two copies, respectively. Thus, we conclude that genetic diversity of the species *V. cholerae* derives predominantly by lateral gene transfer, of which several such transfers appear to be cassettes.

#### **Genomic definition of *V. cholerae* core genome (CG) clade and pandemic strains.**

The *V. cholerae* CG clade, with both 6th and 7th pandemic strains, is defined by gene content as follows. Twenty-seven genes were present exclusively in the genomes of the CG strains, but only five genes are unique to the CG-1 subclade. Four of these (VCA0198-VCA0201) comprise a genomic island (GI-5) on the small chromosome, including genes coding for cytosine-specific DNA methyltransferase<sup>28</sup> and hypothetical proteins, adjacently located to the IS1004 transposase gene. The 7P strains are differentiated in harboring two unique GIs, the *Vibrio* seventh pandemic island-1 (VSP-1) and VSP-2, first discovered by microarray analysis<sup>29</sup>. In addition to the 7P strains, a variant of VSP-1 was found in *V. cholerae* biovar albensis VL426 (Supplementary Fig. 3). Similarly, VSP-2 like GIs were detected in three non-CG strains (TMA21, O39 MZO-3, O135 RC385). Interestingly, a similar GI was also detected in *Vibrio vulnificus* YJ016 and *Vibrio splendidus* 12B01, suggesting that VSP-2 may

be widespread in its distribution among vibrios (Supplementary Fig. 4). It should be noted that stability of these well known pathogenicity islands among 7P members is questionable since most of VPI-2 and VSP-2 were found to be deleted in MO10 and CIRS 101, respectively.

*V. cholerae* contains a super-integron, a large integron island (gene capture system), in the small chromosome (~120 Kbp), a region comprising predominantly hypothetical genes and proposed as a source of genetic variation <sup>10</sup>. All *V. cholerae* strains examined in this study have this integron, a source of much of the variation in gene content (Supplementary Fig. 5). Interestingly, if this region is excluded, all six members of the 7P clade have an identical gene content, with exception of a few genomic islands, including those found in the CTX $\phi$  attachment region. An SXT element, belonging to a family of conjugative transposon-like mobile genetic elements, encodes multiple antibiotic resistance genes and is present only in *V. cholerae* MO10, CIRS 101, MJ-1236, and B33, but not in the other *V. cholerae* strains. *V. cholerae* O139 MO10 differs from other members of the 7P clade in having an O139 antigen specific genomic island, a finding that strongly supports the conclusion of several previous studies, namely that *V. cholerae* O139 derives from a 7th pandemic *V. cholerae* O1 El Tor strain <sup>30</sup>. No other *V. cholerae* O139-specific genes were found in *V. cholerae* MO10.

The "hybrid" strains, possessing an El Tor biotype phenotype, but classical biotype CTX $\phi$ , were isolated during current cholera epidemics in Asia and Africa <sup>6,7</sup>. Two hybrid



strains (B33 and MJ-1236) share a virtually identical genome backbone. Among 3,587,239 bp of orthologous protein-coding region, only 106 nucleotide positions are different and the only significant difference is the presence of a *V. cholerae* MJ-1236 specific 19,729 bp genomic island (GI-12). This GI occurs four times as an almost identical sequence in the large chromosome, with 14 genes including those coding for the putative phage integrase and type I restriction-modification system, implying that it is probably a recently introduced novel temperate bacteriophage. It is not clear why these hybrid strains outcompete other *V. cholerae* O1 El Tor/O139 in the clinical setting, but a key to the puzzle surely lies in the differences among closely related strains, i.e., tandem copies of CTX<sup>Class</sup>, GI-14 and single nucleotide polymorphisms. In addition to these hybrid clones, *V. cholerae* O1 El Tor strains producing a more virulent type of cholera toxin subunit B, i.e., classical, repeatedly have been isolated from patients in Asia and Africa <sup>6</sup>. The genome sequence of a representative of this newly emerged group, i.e. *V. cholerae* CIRS 101, reveals that these strains also have a typical 7P gene content, but with CTX<sup>El Tor</sup>, not CTX<sup>Class</sup>, albeit expressing the classical type subunit B protein (Fig. 3).

The comparative genomics of phylogenetically diverse strains has permitted analysis of the mechanism by which current 7th pandemic clones may have arisen. An highly conserved gene content, synteny, and significant similarity among the six strains of the 7P clade indicate that these *V. cholerae* strains share an almost identical genome "backbone",

having evolved very recently from a common ancestral strain. An hypothetical evolutionary pathway proposed for *V. cholerae* (Fig. 5), with GI migration matched to a genome-based phylogenetic tree, allows the conclusion that the ancestor for the 7P clade was a *V. cholerae* O1 El Tor strain containing several GIs (VPI-1,2, GI-1 to GI-10), receiving VSP-1, VSP-2 and GI-11(=Kappa prophage) by lateral gene transfer, and finally giving rise to the contemporary *V. cholerae* O1 El Tor and O139 strains. It is interesting to note that such an hypothetical ancestral strain shows a gene content similar to *V. cholerae* O1 El Tor BX330286 isolated from a water sample collected in Australia in 1986, a geographic location near Indonesia where the first 7th pandemic *V. cholerae* O1 El Tor was reported in 1961.

#### **Mechanism of *Vibrio cholerae* evolution.**

Of the few human pathogens for which multiple genomes have been fully sequenced<sup>14,12,31</sup>, *V. cholerae* provides a unique opportunity to elucidate evolutionary mechanism(s) that can be generally applied to other bacterial species, since *V. cholerae* is both highly pathogenic for humans and a successful inhabitant of the natural environment worldwide. It is indigenous to estuarine environments of both cholera epidemic and non-epidemic countries<sup>3</sup>.

Unlike *Samonella* Typhi and *Bacillus anthracis*, which are special cases of bacterial species showing clonal properties, *V. cholerae*, with *Streptococcus agalactiae* and *Escherichia coli*, offers a prime example of the important role of lateral gene transfer in the

evolution of a bacterial species. The transition from 6th to 7th cholera pandemic genome type is concluded to result from a change of causal agents of *V. cholerae* O1 classical to *V. cholerae* O1 El Tor biotype. We propose the term "shift" for the event occurring between two distinct phyletic lineages (Fig. 1a). In contrast, the present cholera global pandemic is ascribed to a clonal shift/change among 7P strains, e.g. emergence of *V. cholerae* O139, *V. cholerae* O1 El Tor hybrid and *V. cholerae* O1 El Tor with an altered cholera toxin subunit B. These represent transitions among genetically nearly identical clones, with a few different GIs, for which we propose the term "drift". Much as in the case of the influenza viruses, cholera bacteria undergo a shift/drift cycle over time, though the drift in *V. cholerae* is derived mainly from lateral gene transfer, but not mutation, most likely occurring in the natural environment, where the reservoir of *V. cholerae* occurs in association with its plankton hosts <sup>3,32</sup>.

The present cholera global pandemic is concluded to have been initiated by multiple descendants of a *V. cholerae* O1 El Tor ancestor, diversified and continuously rapidly evolving, mainly *via* lateral gene transfer and most likely driven by environmental factors. Most importantly, the common genome backbone and variable genomic islands of the 7P clade of *V. cholerae* require that a reevaluation be done of the epidemiological practice that employs serogroups as the primary marker for *V. cholerae*. The so-called pandemic clones, instead, should be defined by gene content, the description of which offers significantly

greater potential for the development of reliable and useful diagnostics, vaccines, and therapeutics for cholera. Without doubt, more variants of the 7P clade, as a result of "drift", will be encountered in the future, yielding new serogroups (other than O1 and O139) and phenotypic combinations. Public health workers will be unprepared if the evolution of this species remains unappreciated as an ongoing process in the natural environment, where *V. cholerae* is autochthonous and plays an important role in the nutrient cycles of the natural aquatic ecosystem.

## **METHODS SUMMARY**

Nearly complete genome sequences were obtained from a blend of Sanger, 454 and Solexa sequences using standard protocols used in the DOE Joint Genome Institute. Gene-finding and annotation were achieved using the RAST server<sup>11</sup>. Pairwise genome by genome comparison was carried out using BLASTN, BLASTP and TBLASTX analyses of each ORFs. Additionally, orthologous regions were identified by pairwise global alignment, and used for translation-frame-independent means of comparison. Orthologs and paralogs were differentiated by reciprocal comparison. Sets of orthologous genes were aligned separately by CLUSTALW2 and subsequently concatenated to generate phylogenetic trees using Kimura-2-parameter model<sup>33</sup> and neighbor-joining method<sup>34</sup>. Identification of genomic islands were carried out manually by comparing ortholog content among test strains and phylogenetic tree based on genome sequence.

## **Acknowledgements**

This study was supported by the KOSEF National Research Laboratory Program (Grant No. R0A-2005-000-10110-0 to J.C.), National Institutes of Health (Grant No. 1R01A139129-01 to R.R.C.), National Oceanic and Atmospheric Administration, Oceans and Human Health Initiative (Grant No. S0660009 to R.R.C.), IC Post-Doctoral Fellowship Program (to C.J.G) and the Korean and Swedish governments (to IVI). Funding for sequencing was provided by the Office of the Chief Scientist (USA).

## METHODS

**Genome sequencing.** Draft sequences were obtained from a blend of Sanger and 454 sequences and involved paired end Sanger sequencing on 8kb plasmid libraries to 5X coverage, 20X coverage of 454 data, and optional paired end Sanger sequencing on 35kb fosmid libraries to 1-2X coverage (depending on repeat complexity). To finish the genomes, a collection of custom software and targeted reaction types were used. In addition to targeted sequencing strategies, Solexa data in an untargeted strategy were used to improve low quality regions and to assist gap closure. Repeat resolution was performed using in-house custom software (dupFinisher, C. Han, unpublished). Targeted finishing reactions included transposon bombs<sup>35</sup>, primer walks on clones, primer walks on PCR products, and adapter PCR reactions. Gene-finding and annotation were achieved using the RAST server<sup>11</sup>.

**Comparative genomics.** Genome to genome comparison was performed using three approaches, since completeness and quality of nucleotide sequences varied from strain to strain in the set examined in this study. Firstly, nucleotide sequences as whole contigs were directly aligned using the mummer program<sup>36</sup>. Secondly, ORFs of a given pair of genomes were reciprocally compared each other, using the BLASTN, BLASTP and TBLASTX programs (ORF-dependent comparison). Thirdly, a bioinformatic pipeline was developed to identify homologous regions of a given query ORF. Initially, a segment on target contig homologous to a query ORF was identified using the BLASTN program. This potentially

homologous region was then expanded in both directions by 2,000 bp. Then, nucleotide sequences of the query ORF and selected target homologous region were aligned using a pairwise global alignment algorithm <sup>37</sup>, and the resultant matched region in subject contig was extracted and saved as a homolog (ORF-independent comparison). Orthologs and paralogs were differentiated by reciprocal comparison. In most cases, both ORF-dependent and -independent comparisons yielded the same orthologs, though ORF-independent method performed better for draft sequences of low quality, in which sequencing errors, albeit rare, hampered identification of correct ORFs.

**Identification and annotation of genomic islands.** In this study, we defined genomic islands (GIs) as a continuous array of five or more ORFs that were found to be discontinuously distributed among genomes of test strains. Correct transfer or insertion of GIs was readily differentiated from deletion event by comparing genome-based phylogenetic tree and full matrices showing pairwise detection of orthologous genes between test strains. Identified GIs were designated, and annotated using the BLASTP search of its member ORFs against Genbank NR database.

**Phylogenetic analyses based on genome sequences.** A set of orthologues for each ORF of *V. cholerae* N16961 was obtained for different sets of strains, and then aligned individually using the CLUSTALW2 <sup>38</sup> program. The resultant multiple alignments were concatenated to generate genome scale alignments, which were subsequently used to reconstruct the

neighbor-joining phylogenetic tree<sup>34</sup>. The evolutionary model of Kimura<sup>33</sup> was used to generate the distance matrix. The MEGA program<sup>39</sup> was used for phylogenetic analysis.



## FIGURE LEGENDS

Figure 1. Neighbor-joining trees showing phylogenetic relationships of 23 *V. cholerae* strains representing diverse serogroups. (a) All *V. cholerae* strains based on 1,676 ORFs (1,370,469 bp) (b) core genome (CG) clade based 2,663 ORFs (2,567,393 bp) (c) 7<sup>th</sup> pandemic (7P) clade based on 3,364 ORFs (3,291,577 bp). Bootstrap supports, as percentage, are indicated at the branching points. Bars represent the numbers of substitution per site, respectively. Only orthologous genes showing >95% nucleotide sequence similarity to those of *V. cholerae* N16961 were selected. The tree was rooted using *Vibrio vulnificus* YJ016 and *Vibrio parahaemolyticus* RIMD 2210633.

Figure 2. (a) Neighbor-joining tree of gene cluster for core oligosaccharide (OS) biosynthesis based on 3,982 bp long alignment, which was generated from four ORFs (VC0227, VC0234, VC0236, VC0239) commonly present in all 23 *V. cholerae* strains. (b) Organization of *wav* and *wb\** gene clusters of *V. cholerae*. The homologous regions in the O antigen coding gene cluster are indicated in the same colors. The nomenclature of *wav* gene cluster coding for core oligosaccharide is according to Nesper et al. Strain VL426 has a variant of type 3 which we designate type 3b in this study. (c) Neighbor-joining tree based on 19 genes (VC0241-VC0254, VC0259-VC0263) coding for O1 antigen. The numbers at nodes indicate the bootstrap supports

and bars indicate the number of substitutions per a nucleotide position.

Figure 3. Schematic representation of various prophages and genetic elements present in the target regions of CTX $\phi$  insetion. The numbers in parentheses are year of isolation and source (C, clinical; E, Environment). \* TLC, El Tor type CTX  $\phi$ , RS1 element are found, but no positional information can be obtained from assemblies. † Classical type CTX  $\phi$  and RS1 are present, but no positional information can be obtained.

Figure 4 . Genomic representation of genomic islands of both *V. cholerae* chromosomes. The two circles in the middle represent the genes in *V. cholerae* O1 El Tor N16961. The inner circle indicates genomic islands found in strain N16961, whereas the outer circles are those absent in strain N16961.

Figure 5. Proposed hypothetical evolutionary pathway of the *V. cholerae* species. Probable insertions and deletions of genomic islands (Supplementary Table 2) found in 23 *V. cholerae* strains are indicated by black and red arrows, respectively, along the phylogenetic tree based on genome sequence data. Hypothetical ancestral strains are indicated by open circles.

## REFERENCES

- <sup>1</sup> Chatterjee, S. N. & Chaudhuri, K. Lipopolysaccharides of *Vibrio cholerae*.  
I. Physical and chemical characterization. *Biochim Biophys Acta* **1639**,  
65-79 (2003).
- <sup>2</sup> Kaper, J. B., Morris, J. G., Jr. & Levine, M. M. Cholera. *Clin Microbiol  
Rev* **8**, 48-86 (1995).
- <sup>3</sup> Colwell, R. R. Global climate and infectious disease: the cholera  
paradigm. *Science* **274**, 2025-2031 (1996).
- <sup>4</sup> Ramamurthy, T. *et al.* Emergence of novel strain of *Vibrio cholerae* with  
epidemic potential in southern and eastern India. *Lancet* **341**, 703-704  
(1993).
- <sup>5</sup> Nair, G. B. *et al.* New variants of *Vibrio cholerae* O1 biotype El Tor with  
attributes of the classical biotype from hospitalized patients with acute  
diarrhea in Bangladesh. *J Clin Microbiol* **40**, 3296-3299 (2002).
- <sup>6</sup> Nair, G. B. *et al.* Cholera due to altered El Tor strains of *Vibrio cholerae*  
O1 in Bangladesh. *J Clin Microbiol* **44**, 4211-4213, doi:JCM.01304-06 [pii]  
10.1128/JCM.01304-06 (2006).
- <sup>7</sup> Ansaruzzaman, M. *et al.* Cholera in Mozambique, variant of *Vibrio*

- cholerae. *Emerg Infect Dis* **10**, 2057-2059 (2004).
- <sup>8</sup> Waldor, M. K. & Mekalanos, J. J. Lysogenic conversion by a filamentous phage encoding cholera toxin. *Science* **272**, 1910-1914 (1996).
- <sup>9</sup> Blokesch, M. & Schoolnik, G. K. Serogroup conversion of *Vibrio cholerae* in aquatic reservoirs. *PLoS Pathog* **3**, e81 (2007).
- <sup>10</sup> Heidelberg, J. F. *et al.* DNA sequence of both chromosomes of the cholera pathogen *Vibrio cholerae*. *Nature* **406**, 477-483 (2000).
- <sup>11</sup> Aziz, R. K. *et al.* The RAST Server: rapid annotations using subsystems technology. *BMC Genomics* **9**, 75 (2008).
- <sup>12</sup> Rasko, D. A. *et al.* The pangenome structure of *Escherichia coli*: comparative genomic analysis of *E. coli* commensal and pathogenic isolates. *J Bacteriol* **190**, 6881-6893, doi:JB.00619-08 [pii]10.1128/JB.00619-08 (2008).
- <sup>13</sup> Chen, S. L. *et al.* Identification of genes subject to positive selection in uropathogenic strains of *Escherichia coli*: a comparative genomics approach. *Proc Natl Acad Sci U S A* **103**, 5977-5982, doi:0600938103 [pii]10.1073/pnas.0600938103 (2006).
- <sup>14</sup> Tettelin, H. *et al.* Genome analysis of multiple pathogenic isolates of *Streptococcus agalactiae*: implications for the microbial "pan-genome".

*Proc Natl Acad Sci U S A* **102**, 13950-13955, doi:0506758102

[pii]10.1073/pnas.0506758102 (2005).

- <sup>15</sup> Nesper, J. *et al.* Comparative and genetic analyses of the putative *Vibrio cholerae* lipopolysaccharide core oligosaccharide biosynthesis (*wav*) gene cluster. *Infect Immun* **70**, 2419-2433 (2002).
- <sup>16</sup> Li, M., Shimada, T., Morris, J. G., Jr., Sulakvelidze, A. & Sozhamannan, S. Evidence for the emergence of non-O1 and non-O139 *Vibrio cholerae* strains with pathogenic potential by exchange of O-antigen biosynthesis regions. *Infect Immun* **70**, 2441-2453 (2002).
- <sup>17</sup> Yamasaki, S. *et al.* The genes responsible for O-antigen synthesis of *vibrio cholerae* O139 are closely related to those of *vibrio cholerae* O22. *Gene* **237**, 321-332, doi:S0378111999003443 [pii] (1999).
- <sup>18</sup> Safa, A. *et al.* Multilocus genetic analysis reveals that the Australian strains of *Vibrio cholerae* O1 are similar to the pre-seventh pandemic strains of the El Tor biotype. *J. Med. Microbiol.* **58**, 105-111 (2009).
- <sup>19</sup> Farfan, M., Minana, D., Fuste, M. C. & Loren, J. G. Genetic relationships between clinical and environmental *Vibrio cholerae* isolates based on multilocus enzyme electrophoresis. *Microbiology* **146** ( Pt 10), 2613-2626 (2000).

- <sup>20</sup> Bik, E. M., Gouw, R. D. & Mooi, F. R. DNA fingerprinting of *Vibrio cholerae* strains with a novel insertion sequence element: a tool to identify epidemic strains. *J Clin Microbiol* **34**, 1453-1461 (1996).
- <sup>21</sup> Meibom, K. L., Blokesch, M., Dolganov, N. A., Wu, C. Y. & Schoolnik, G. K. Chitin induces natural competence in *Vibrio cholerae*. *Science* **310**, 1824-1827 (2005).
- <sup>22</sup> Colwell, R. R., Huq, A., Chowdhury, M. A., Brayton, P. R. & Xu, B. Serogroup conversion of *Vibrio cholerae*. *Can J Microbiol* **41**, 946-950 (1995).
- <sup>23</sup> Davis, B. M., Moyer, K. E., Boyd, E. F. & Waldor, M. K. CTX prophages in classical biotype *Vibrio cholerae*: functional phage genes but dysfunctional phage genomes. *J Bacteriol* **182**, 6992-6998 (2000).
- <sup>24</sup> Udden, S. M. *et al.* Acquisition of classical CTX prophage from *Vibrio cholerae* O141 by El Tor strains aided by lytic phages and chitin-induced competence. *Proc Natl Acad Sci U S A* **105**, 11951-11956, doi:0805560105 [pii] 10.1073/pnas.0805560105 (2008).
- <sup>25</sup> Rubin, E. J., Lin, W., Mekalanos, J. J. & Waldor, M. K. Replication and integration of a *Vibrio cholerae* cryptic plasmid linked to the CTX

- prophage. *Mol Microbiol* **28**, 1247-1254 (1998).
- <sup>26</sup> Faruque, S. M. *et al.* Genomic analysis of the Mozambique strain of *Vibrio cholerae* O1 reveals the origin of El Tor strains carrying classical CTX prophage. *Proc Natl Acad Sci U S A* **104**, 5151-5156 (2007).
- <sup>27</sup> Faruque, S. M. *et al.* CTXphi-independent production of the RS1 satellite phage by *Vibrio cholerae*. *Proc Natl Acad Sci U S A* **100**, 1280-1285, doi:10.1073/pnas.02373851000237385100 [pii] (2003).
- <sup>28</sup> Banerjee, S. & Chowdhury, R. An orphan DNA (cytosine-5-)-methyltransferase in *Vibrio cholerae*. *Microbiology* **152**, 1055-1062, doi:152/4/1055 [pii] 10.1099/mic.0.28624-0 (2006).
- <sup>29</sup> Dziejman, M. *et al.* Genomic characterization of non-O1, non-O139 *Vibrio cholerae* reveals genes for a type III secretion system. *Proc Natl Acad Sci U S A* **102**, 3465-3470 (2005).
- <sup>30</sup> Karaolis, D. K., Lan, R. & Reeves, P. R. Molecular evolution of the seventh-pandemic clone of *Vibrio cholerae* and its relationship to other pandemic and epidemic *V. cholerae* isolates. *J Bacteriol* **176**, 6199-6206 (1994).
- <sup>31</sup> Holt, K. E. *et al.* High-throughput sequencing provides insights into

- genome variation and evolution in Salmonella Typhi. *Nat Genet* **40**, 987-993, doi:ng.195 [pii]10.1038/ng.195 (2008).
- <sup>32</sup> Constantin de Magny, G. *et al.* Environmental signatures associated with cholera epidemics. *Proc Natl Acad Sci U S A* **105**, 19676-17681 (2008).
- <sup>33</sup> Kimura, M. A simple method for estimating evolutionary rate of base substitutions through comparative studies of nucleotide sequences. *J Mol Evol* **16**, 111-120 (1980).
- <sup>34</sup> Saitou, N. & Nei, M. The neighbor-joining method: a new method for reconstructing phylogenetic trees. *Mol Biol Evol* **4**, 406-425 (1987).
- <sup>35</sup> Goryshin, I. Y. & Reznikoff, W. S. Tn5 in vitro transposition. *J Biol Chem* **273**, 7367-7374 (1998).
- <sup>36</sup> Kurtz, S. *et al.* Versatile and open software for comparing large genomes. *Genome Biol* **5**, R12 (2004).
- <sup>37</sup> Myers, E. W. & Miller, W. Optimal alignments in linear space. *Comput Appl Biosci* **4**, 11-17 (1988).
- <sup>38</sup> Larkin, M. A. *et al.* Clustal W and Clustal X version 2.0. *Bioinformatics* **23**, 2947-2948 (2007).
- <sup>39</sup> Kumar, S., Nei, M., Dudley, J. & Tamura, K. MEGA: a biologist-centric software for evolutionary analysis of DNA and protein sequences. *Brief*



*Bioinform* **9**, 299-306 (2008).

**Supplementary Information** is linked to the online version of the paper at

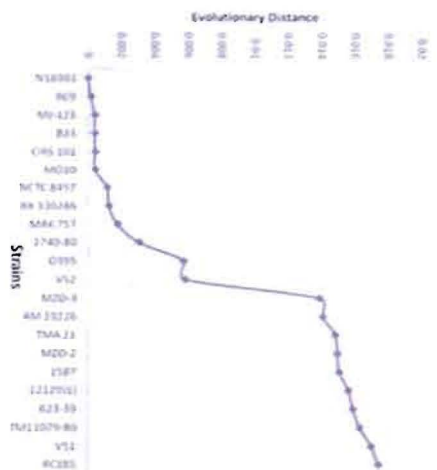
[www.nature.com/nature](http://www.nature.com/nature).



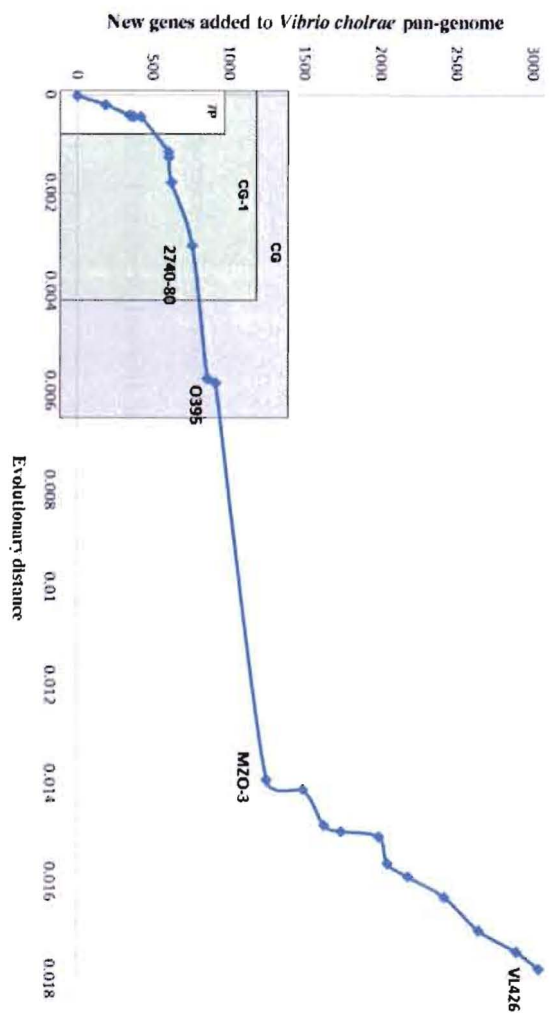


[illegible]



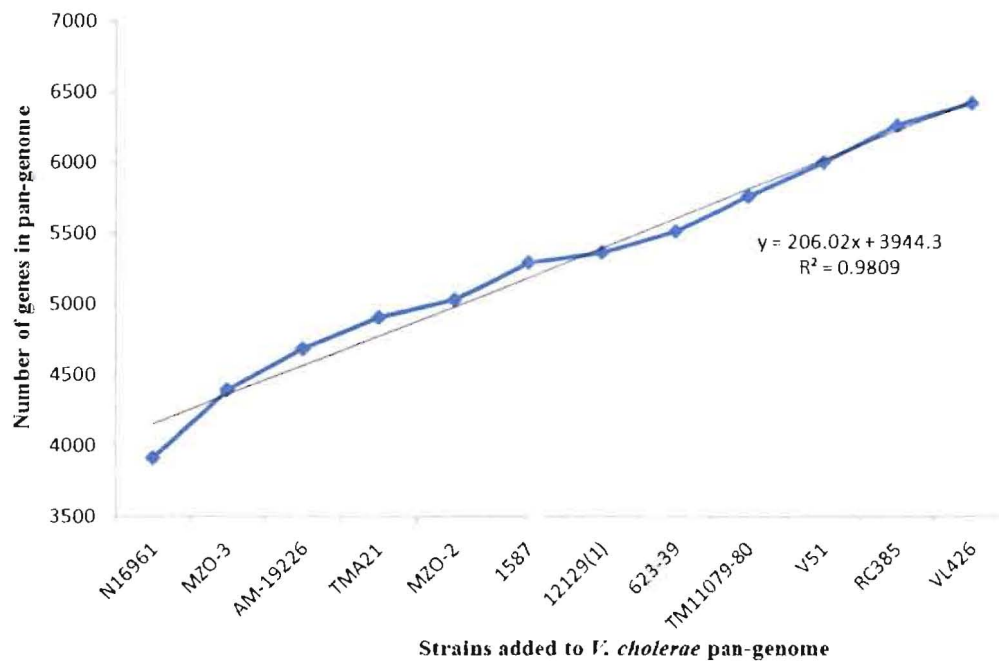


(ii)



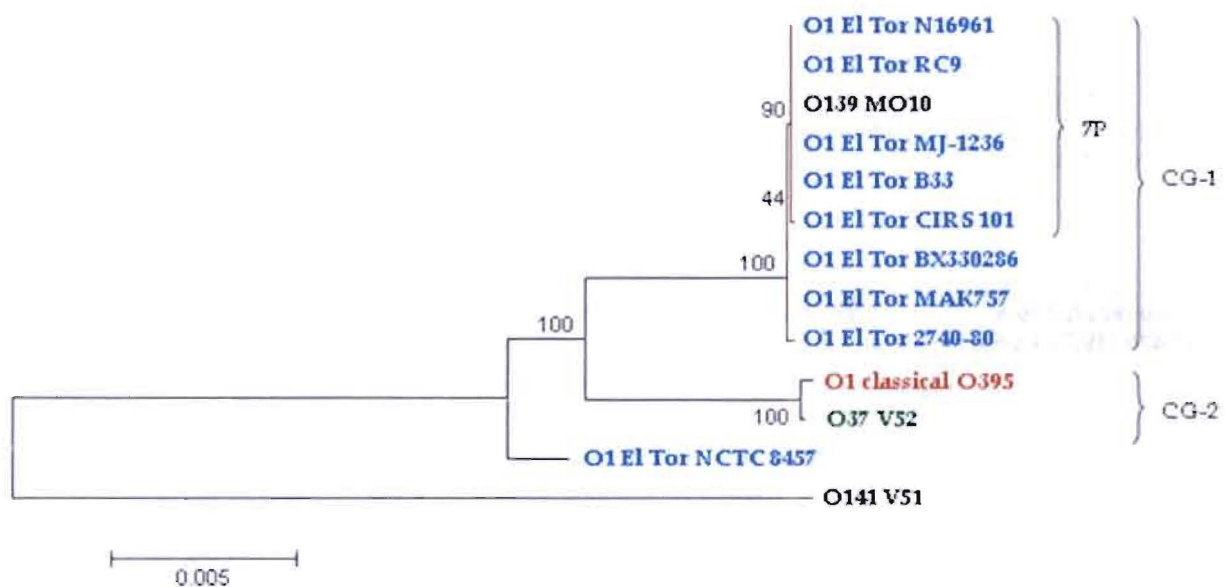
(b)



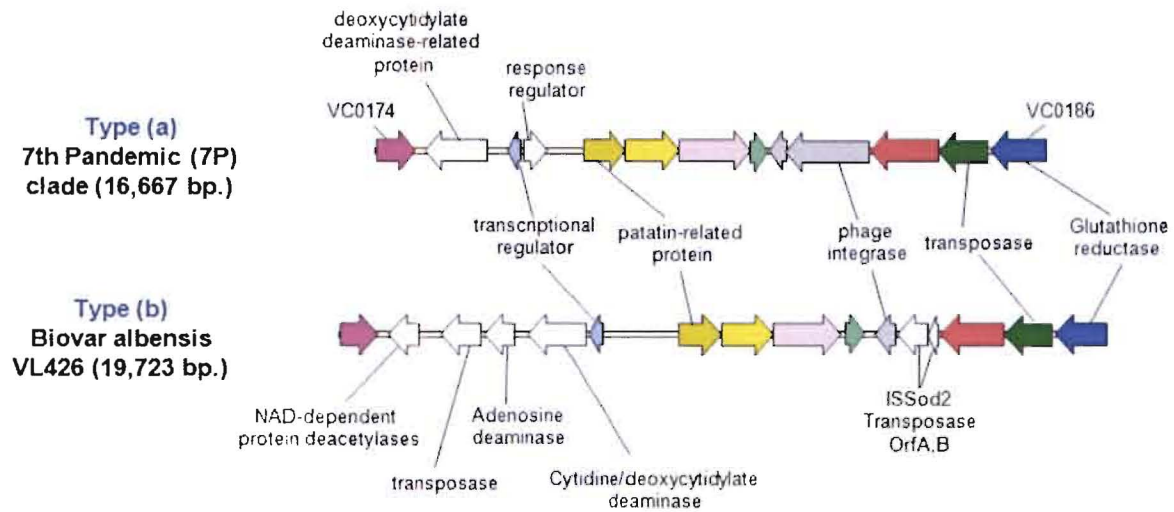


(c)

**Supplementary Figure 1.** Pan-genome analysis of *Vibrio cholerae*. (a) Evolutionary distances of *V. cholerae* strains from the reference strain (O1 El Tor N16961). (b) cumulative number of newly found additional genes to *V. cholerae* pan-genome with N16961 as reference. Strains were added in the order as shown in (a). (c) Pan-genome calculated from strains that are representatives of 12 distinct phyletic lines (Fig. 1a). Linear regression indicates that 206 genes will be added to the pan-genome of *V. cholerae*, if new strain represents equally distinct phyletic lineage.



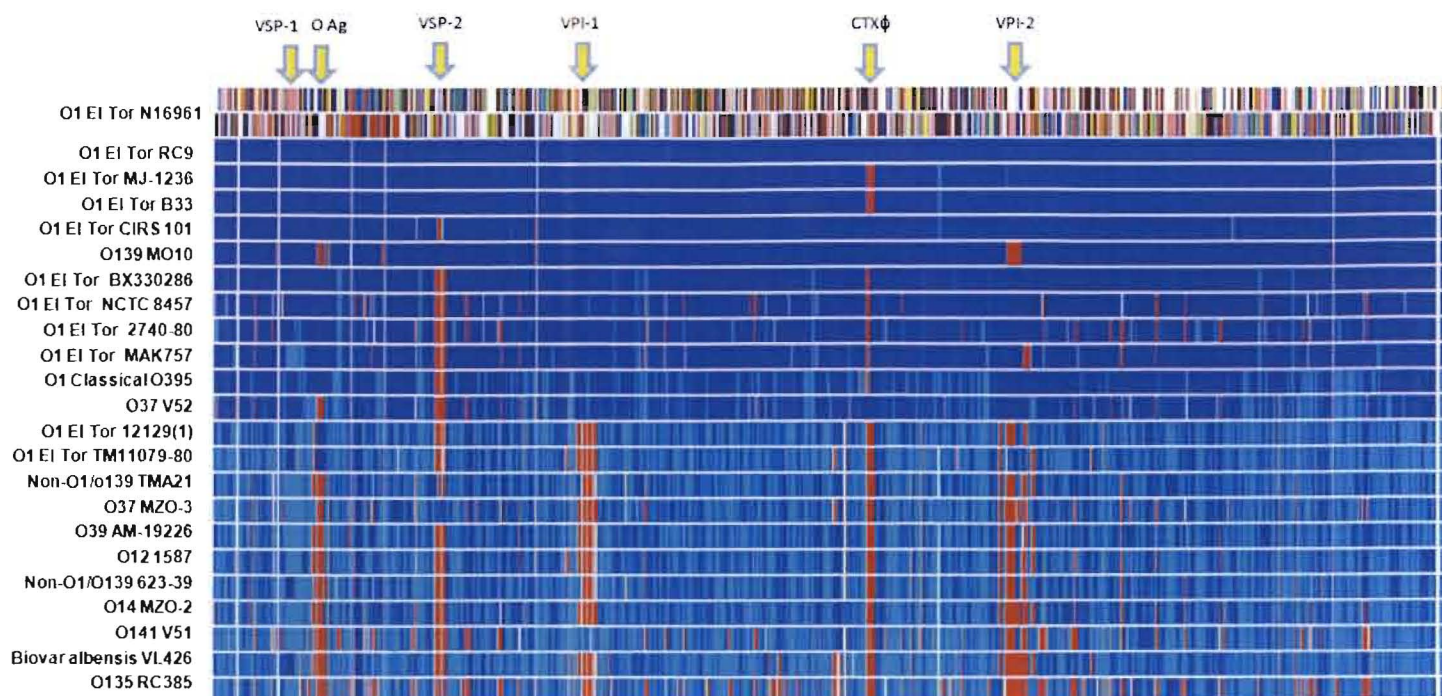
**Supplementary Figure 2.** Neighbor-joining tree showing evolutionary relationships of *Vibrio* pathogenicity island (VPI-1) found in *V. cholerae* strains. The calculation was based on 24 orthologous ORFs (VC0821-VC0834, VC0836-VC0845) comprising 27,393 bp. The phylogenetic placement of *V. cholerae* O1 El Tor NCTC 8457 (CG-1 member) indicates the replacement of its VPI-1 by that from a non-CG strain. Bar represents 0.005 substitution per site.



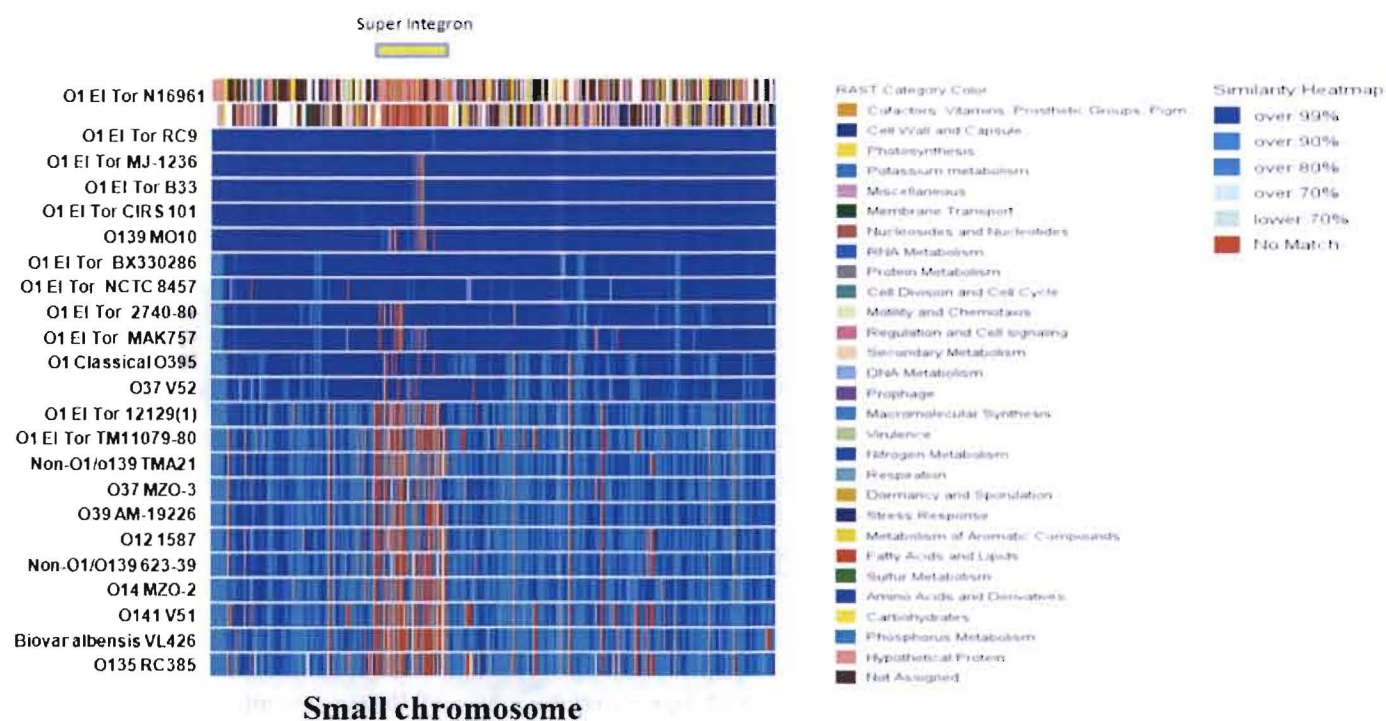
**Supplementary Figure 3.** Genetic organization of *Vibrio* seventh pathogenicity island-1 (VSP-1) genomic island (N16961 ORFs: VC0175-VC0185) found in only 7th pandemic (7P) *V. cholerae* strains and *V. cholerae* biovar albensis VL426.





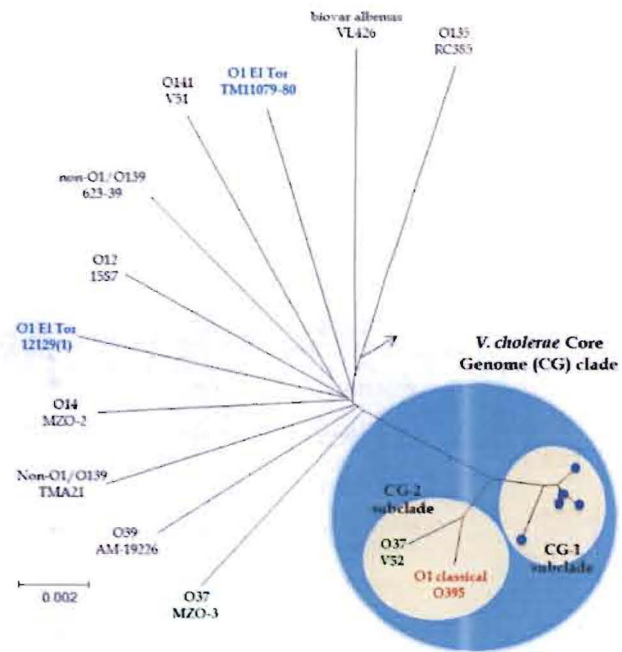


## Large chromosome

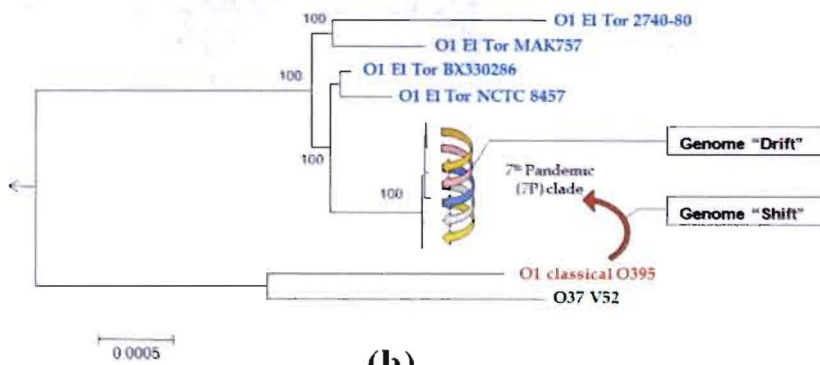


## Small chromosome

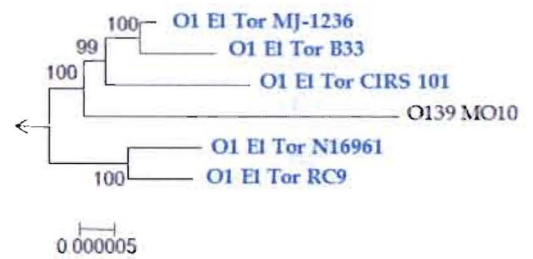
**Supplementary Figure 5.** Gene contents of 22 *V. cholerae* strains showing gene conservation using *V. cholerae* O1 El Tor N16961 as the reference strain. Orthologues found are color-coded according to nucleotide sequence similarity values to the reference ORFs. Major genomic islands are indicated at the top row.



(a)



(b)



(c)

Figure 1.



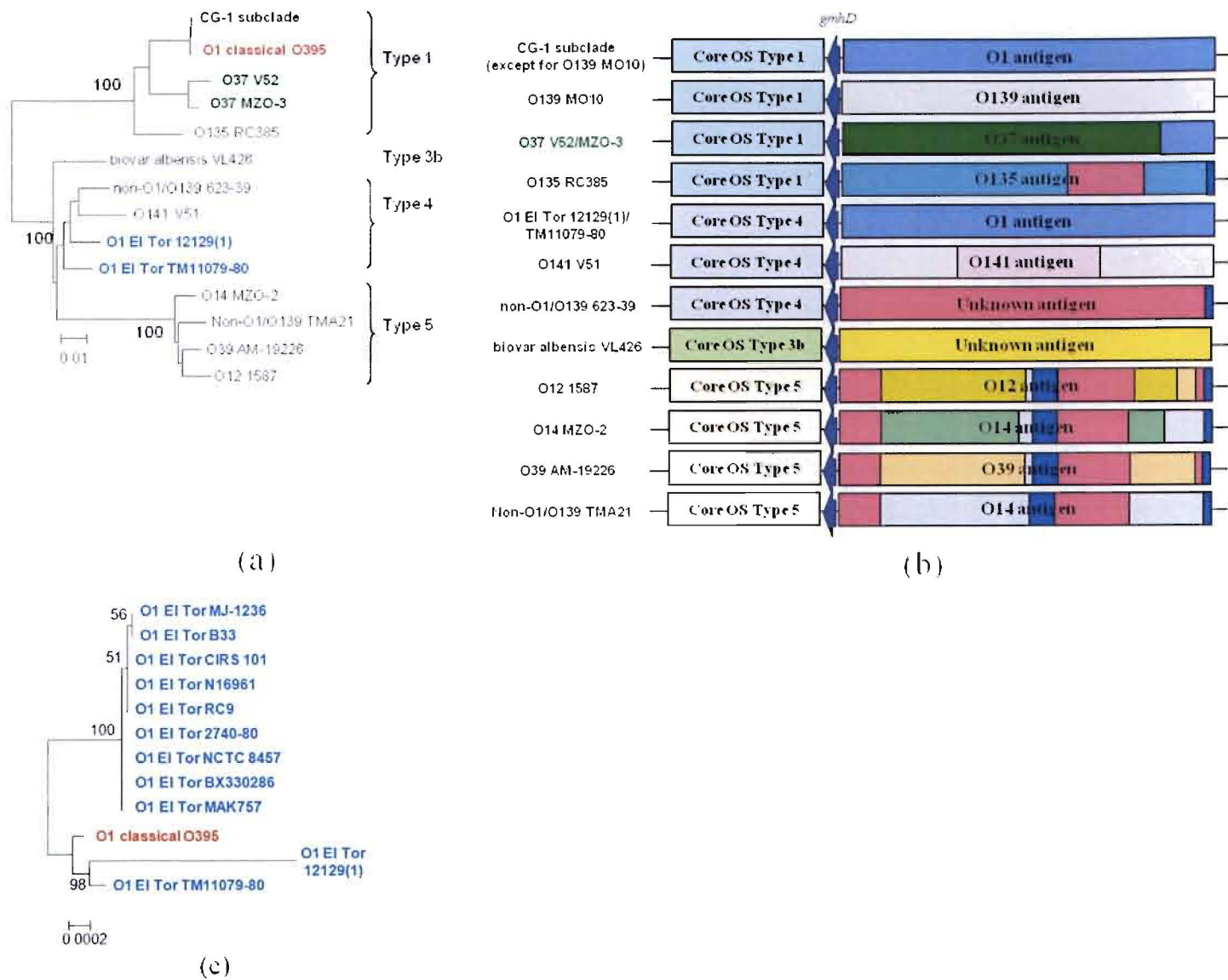


Figure 2.

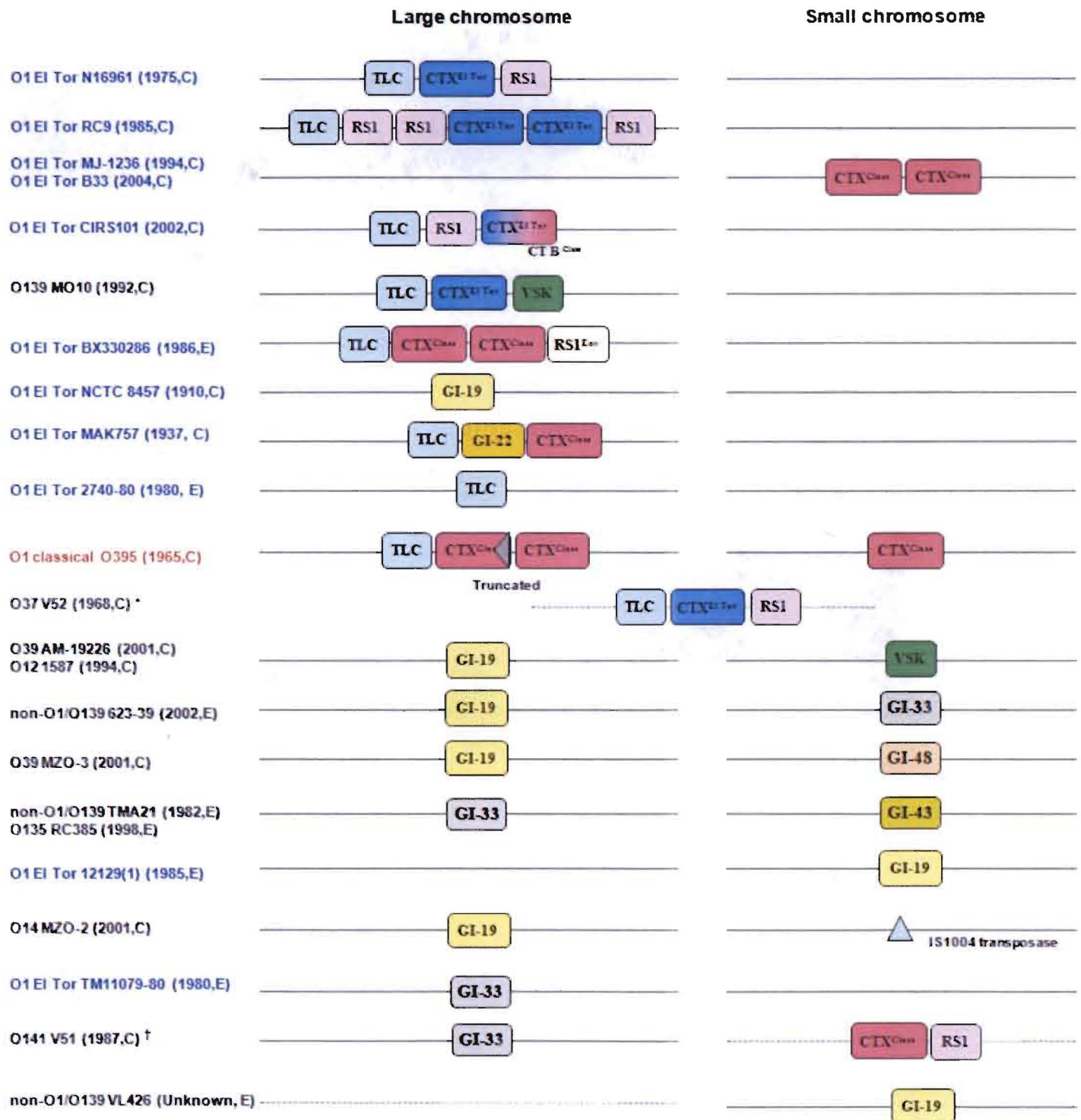


Figure 3.



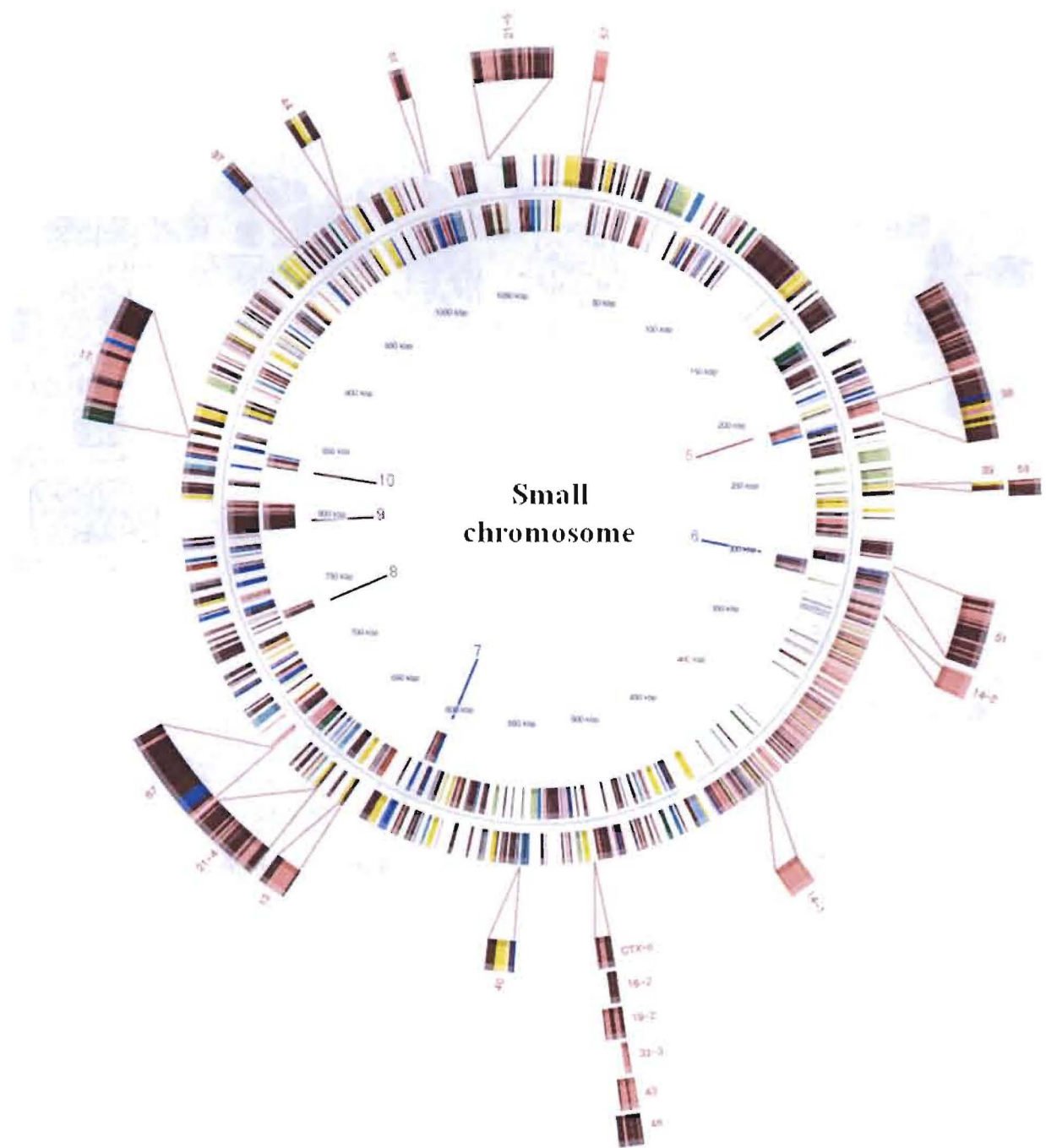


Figure 4 .



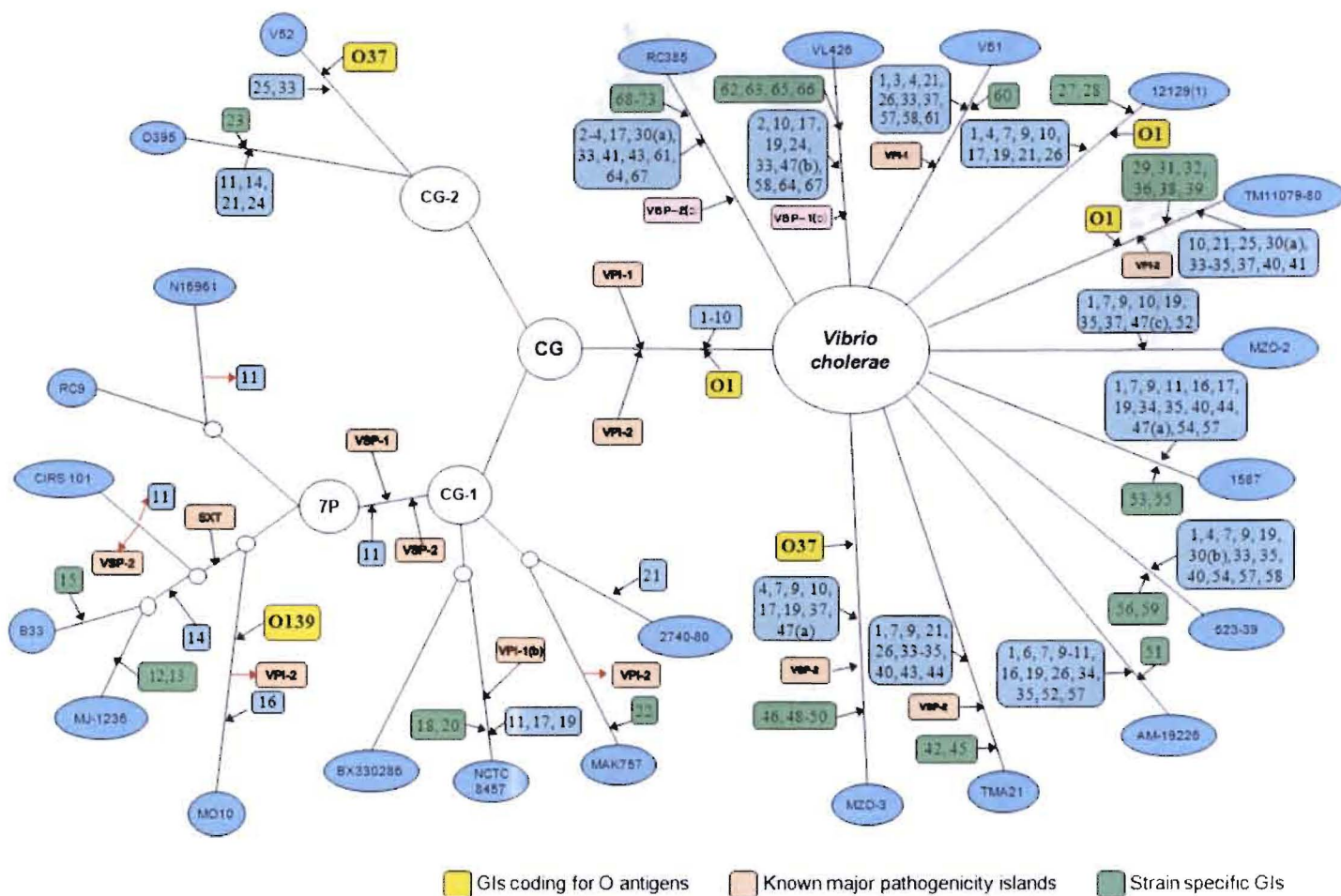


Figure 5.



**Table 1. Characteristics of *Vibrio cholerae* strains compared in this study.**

Strain	Genome Code	Serogroup	Biotype	Geographical Origin	Source of isolation	Year of isolation	Sequencing status <sup>a</sup>	Number of contigs	Accession
N16961	VCN16961	O1 Inaba	El Tor	Bangladesh	Clinical	1975	Complete	2	NC_002505/ NC_002506
RC9	VCRC9	O1 Ogawa	El Tor	Kenya	Clinical	1985	S/4/E	11	This study
MJ-1236	VCMJ1236	O1 Inaba	El Tor	Matlab, Bangladesh	Clinical	1994	Complete	2	This study
B33	VCB33	O1 Ogawa	El Tor	Beira, Mozambique	Clinical	2004	S/4/E	17	This study
CIRS 101	VCCIRS101	O1 Inaba	El Tor	Dhaka, Bangladesh	Clinical	2002	S/4/E	18	This study
MO10	VCMO10	O139		Madras, India	Clinical	1992	S/4	84	AAKF03000000
2740-80	VC274080	O1 Inaba	El Tor	US Gulf Coast	Water	1980	Sanger	257	AAUT01000000
BX330286	VCBX330286	O1 Inaba	El Tor	Australia	Water	1986	Complete	8	This study
MAK757	VCMAK757	O1 Ogawa	El Tor	Celebes Islands	Clinical	1937	Sanger	206	AAUS00000000
NCTC 8457	VC8457	O1 Inaba	El Tor	Saudi Arabia	Clinical	1910	Sanger	390	AAWD01000000
O395	VCO395	O1 Ogawa	Classical	India	Clinical	1965	Complete	2	NC_009456/ NC_009457
V52	VCV52	O37		Sudan	Clinical	1968	Sanger	268	AAKJ02000000
12129(1)	VC12129	O1 Inaba	El Tor	Australia	Water	1985	S/4/E	12	This study
TM 11079-80	VCTM11079	O1 Ogawa	El Tor	Brazil	Sewage	1980	S/4/E	35	This study
VL426	VCVL426	non-O1/O139	albensis	Maidstone, Kent, UK	Water	Unknown	Complete	5	This study
TMA21	VCTMA21	non-O1/O139		Brazil	Seawater	1982	S/4/E	20	This study
1587	VC1587	O12		Lima, Peru	Clinical	1994	Sanger	254	AAUR01000000
RC385	VCRC385	O135		Chesapeake Bay	Plankton	1998	Sanger	550	AAKH02000000
MZO-2	VCMZO2	O14		Bangladesh	Clinical	2001	Sanger	162	AAWF01000000
V51	VCV51	O141		USA	Clinical	1987	Sanger	360	AAKI02000000
MZO-3	VCMZO3	O37		Bangladesh	Clinical	2001	Sanger	292	AAUU01000000
AM-19226	VCAM19226	O39		Bangladesh	Clinical	2001	Sanger	154	AATY01000000
623-39	VC62339	non-O1/O139		Bangladesh	Water	2002	Sanger	314	AAWG00000000

ger, Draft assemblies by Sanger sequencing; S/4, Sanger sequencing and 454 pyrosequencing were combined; , S/4 followed by quality improvement by standard genome sequencing procedures.

**Supplementary Table 3.** *Vibrio cholerae* genomic islands that occupy the same position/area, showing a cassette like property.

	Insertion region	GI-designation	Strains
Large chromosome	VC0174-VC0186	VSP-1	N16961, RC9, MJ-1236, B33, MO10
		GI-50	MZO-3
	VC0489 -VC0517	VSP-2	N16961, RC9, B33, MJ1236, MO10, TMA21, MZO-3
		GI-56	623-39
	VC0809 - VC0848	VPI-1	N16961, RC9, B33, MJ1236, MO10, BX330286, NCTC 8457, 2740-80, MAK757, O395, V52, V51
		GI-47	MZO-3, 1587, MZO-2, VL426
	VC1757 - VC1810	VPI-2	N16961, RC9, MJ1236, B33, MO10, BX330286, NCTC8457, 2740-80, MAK757, O395, V52, TM11079-80
		GI-26	12129(1), TMA21, AM-19226, V51
		GI-46	MZO-3
		GI-54	1587, 623-39
	VC1451 - VC1465	CTX□	N16961, RC9, O395
		GI-16	MO10
		GI-19	NCTC 8457, AM-19226, 1587, 623-39, MZO-2, MZO-3
		GI-22	MAK757
		GI-33	TM11079-80, TMA21, V51, RC385
	VC0002 - VC0003	GI-15	B33
		GI-30	TM11079-80, 623-39, RC385
	VC0289 - VC0290	GI-24	O395, VL426
		GI-42	TMA21
	VC0209 - VC0208	GI-32	TM11079-80
		GI-52	AM-19226, MZO-2
		GI-68	RC385
	VC0847 - VC0807	GI-27	12129(1)
		GI-41	TM11079-80
		GI-59	623-39
	VC1407-VC1414	GI-62	VL426
		GI-72	RC385
Small chromosome	VCA0569 - VCA0570	CTX□	MJ-1236, B33, O395
		GI-16	AM-19226, 1587
		GI-19	12129(1), VL426
		GI-33	623-39
		GI-43	TMA21, RC385
		GI-48	MZO-3
	VCA0197 - VCA0204	GI-5	N16961, RC9, B33, MJ1236, MO10, BX330286, NCTC8457, 2740-80, MAK757, O395, V52
		GI-38	TM11079-80
	VCA0235 - VCA0237	GI-39	TM11079-80
		GI-58	623-39, V51, VL426