

Comparative genomic analysis of 142 bacteriophages infecting *Salmonella enterica* subsp. *enterica*

Ruimin Gao

Canadian Food Inspection Agency

Sohail Naushad

Canadian Food Inspection Agency

Sylvain Moineau

Universite Laval, Quebec City

Lawrence Goodridge

University of Guelph, Guelph

Dele Ogunremi (✉ dele.ogunremi@canada.ca)

Canadian Food Inspection Agency <https://orcid.org/0000-0001-9123-5574>

Research article

Keywords: Comparative genomics, Bacteriophage, Nucleotide identity, *Salmonella enterica*, Phamerator, Prophage sequence typing, Phage clusters

Posted Date: October 11th, 2019

DOI: <https://doi.org/10.21203/rs.2.15923/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License. [Read Full License](#)

Version of Record: A version of this preprint was published at BMC Genomics on May 26th, 2020. See the published version at <https://doi.org/10.1186/s12864-020-6765-z>.

Abstract

Background: Bacteriophages are bacterial parasites and are considered the most abundant and diverse biological entities on the planet. Previously we identified 154 prophages from 151 serovars of *Salmonella enterica* subsp. *enterica*. A detailed analysis of *Salmonella* prophage genomics is required given the influence of phages on their bacterial hosts and should provide a broader understanding of *Salmonella* biology and virulence and contribute to the practical applications of phages as vectors and antibacterial agents.

Results: Comparative analysis of the full genome sequences of 142 prophages of *Salmonella enterica* subsp. *enterica* retrieved from public databases revealed an extensive variation in genome sizes (6.4- 358.7 kb) and guanine plus cytosine (GC) content (35.5-65.4%) and a linear correlation between the genome size and the number of open reading frames (ORFs). We used three approaches to compare the phage genomes. The NUCmer/MUMmer genome alignment tool was used to evaluate linkages and correlations based on nucleotide identity between genomes. Multiple sequence alignment was performed to calculate genome average nucleotide identity using the Kalign program. Finally, genome synteny was explored using dot plot analysis. We found that 90 phage genome sequences grouped into 17 distinct clusters while the remaining 52 genomes showed no close relationships with the other phage genomes and are identified as singletons. We generated genome maps using nucleotide and amino acid sequences which allowed protein-coding genes to be sorted into phamilies (phams) using the Phamerator software. Out of 5796 total assigned phamilies, one phamily was observed to be dominant and was found in 49 prophages, or 34.5% of the 142 phages in our collection. A majority of the phamilies, 4330 out of 5796 (74.7%), occurred in just one prophage underscoring the high degree of diversity among *Salmonella* bacteriophages.

Conclusions: Based on nucleotide and amino acid sequences, a high diversity was found among *Salmonella* bacteriophages which validate the use of prophage sequence analysis as a highly discriminatory subtyping tool for *Salmonella*. Thorough understanding of the conservation and variation of prophage genomic characteristics will facilitate their rational design and use as tools for bacterial strain construction, vector development and as anti-bacterial agents.

Background

The Gram-negative bacterial genus *Salmonella* belongs to the family Enterobacteriaceae, order Enterobacteriales, class Gammaproteobacteria and phylum Proteobacteria. *Salmonella* cells have a length of 2 to 5 μm and a diameter ranging from 0.7 to 1.5 μm , as well as being predominantly motile due to peritrichous flagella [1]. The genus consists of two species, namely *Salmonella enterica* and *S. bongori*. The former can be further divided into six subspecies which corresponds to known serotypes (depicted with Roman numerals): *enterica* (I), *salamae* (II), *arizonae* (IIIa), *diarizonae* (IIIb), *houtenae* (IV) and *indica* (VI) [2]. The serotype V is now considered a separate species and designated *S. bongori*. Based on the presence of somatic O (lipopolysaccharide) and flagellar H antigens (Kauffman-White classification), the above six *S. enterica* subspecies are divided into over 2600 serovars [3] but fewer than 100 serovars have been associated with human illnesses [4]. *Salmonella enterica* subspecies *enterica* is typically categorized into typhoidal and non-typhoidal *Salmonella* as a result of symptoms presenting in infected humans. Non-typhoidal *Salmonella*, which is made up of a large number of the serovars, can be transmitted from animals to humans and between humans, often via vehicles such as foods, and they usually invade only the gastrointestinal tract leading to symptoms that resolve even in the absence of antibacterial therapy [5]. In contrast, typhoidal *Salmonella* serovars such as Typhi, Paratyphi A and Paratyphic C, are transferred from human to human and can cause severe infections requiring antibiotic treatment [6]. Wide spread resistance against antibiotics has prompted a renewed surge of interest in bacteriophages which are viruses capable of infecting and sometimes killing bacteria, as safe and effective therapy alternatives [7].

Bacteriophages, sometimes simply referred to as phages, are considered the most abundant biological entities on the planet [8]. These bacterial viruses can undergo two life cycles: lysis or lysogeny. A bacteriophage capable of only lytic growth is described as virulent. In contrast, temperate bacteriophage refers to the ability of some phages to display a lysogenic cycle and instead of killing the host bacterium becomes integrated into the chromosome. A bacterium that contains a complete set of phage genes is called a lysogen, while the integrated viral DNA is called a prophage. Most temperate phages form lysogens by integration at a unique attachment site in the host chromosome [9, 10]. The integration process has been described as a biological arms race between the infecting virus and the host bacterium [11]. There is an array of host defense mechanisms that are stacked against the virus which in turn increasingly acquires and displays a counter-offensive to thwart and evade the anti-viral mechanisms resulting in integration into the host genome [11–13].

Tailed phages which belong to the Order Caudovirales are the most abundant group of viruses infecting bacteria and are also the most prevalent in the human gut. They are easily recognized under an electron microscope by their polyhedral capsids and tubular tails [14]. The order Caudovirales is made up of five families, namely: (1) *Myoviridae* (contractile tails, long and relatively thick), (2) *Siphoviridae* (long noncontractile tails), (3) *Podoviridae* (short noncontractile tails) [14], (4) *Ackermannviridae* (contractile tails) and (5) *Herelleviridae* - spouna-like (contractile tails, long and relatively thick) [15]. Bacteriophages were first described by Frederick Twort in 1915 and Felix d'Herelle in 1917 [16], and studies into their relationship with *Salmonella enterica* serovar Typhimurium led to the description of "symbiotic bacteriophages" by Boyd [17]. We recently analyzed the bacteriophages present in 1,760 genomes of *Salmonella* strains present in a research database (<https://salfos.ibis.ulaval.ca/>) and apart from three strains devoid of any prophage, the genomes had 1 - 15 prophages with an average of 5 prophages per isolate [18]. Previous analyses of *Salmonella* phages have led to their classification into five groups (P27-like, P2-like, lambdaoid, P22-like, and T7-like) and three outliers (ϵ 15, KS7, and Felix O1) [10]. Apart from the primary role of phage gene products to ensure that these viruses can infect bacteria, survive and reproduce in their hosts, phage genes have been shown to code for virulence factors, toxin, and antimicrobial resistance genes. The presence of these genes appears to contribute in a substantial manner to the evolution of the bacterial host [18–20]. Studies of prophage biology have practical significance in choice of phages as antibacterial agents, in bacterial strain construction and typing for epidemiological purposes [21, 22].

The advent of whole genome sequencing has greatly facilitated the detection and characterization of phages and prophages in bacterial hosts and the ability to evaluate their impacts on the host. Evolutionary analysis of phage genes open reading frames (ORF) families based on sequence analysis of a large number of phage genomes in the GenBank (about 13,703 phage genomes were present as of June 2019) (<http://millardlab.org/bioinformatics/bacteriophage-genomes/phage-genomes-june-2019/>) has provided insights into the impact on the evolution of both the virus and host [23]. Whole-genome comparative

analysis has been successfully applied to study phages present or infecting several bacterial genera including *Mycobacteria* [24], *Staphylococcus* [25], *Bacillus* [26], *Gordonia* [27], *Pseudomonas* [23] and as well as the *Enterobacteriaceae* family [28]. Phage genomes are commonly grouped into clusters, but outlier phages lacking strong nucleotide identity relationships with other clustered genome are often designed as 'singletons' [27]. To classify phage genomes into clusters and subclusters, there are several commonly used tools/approaches. The dot plot program Genome Pair Rapid Dotter (Gepard) [29] can reveal very substantial synteny among genomes. Typically, the dot plot can recognize similarities spanning more than half of the genome lengths [24]. The average nucleotide identity (ANI) are determined using tools such as Kalign [30] and MUMmer [31] using genomes alignment and comparison. Genome map and gene content analyses can be performed using Phamerator, which sorts protein-coding genes into Phamilies (Phams) and generate a database of gene relationships [32, 33].

Using PHASTER (PHAge Search Tool Enhanced Release) [34, 35], we previously demonstrated the presence of 154 different prophages in 1760 *S. enterica* genomes which covered 151 *Salmonella* serovars [18]. We also previously showed that some prophage sequences were conserved among strains belonging to the same serovars and that the prophage repertoires provided an additional marker for differentiating *S. enterica* subtypes during foodborne outbreaks [18]. Here, a more detailed characterization of these *Salmonella* phage genomes was carried out to generate knowledge on their biological variation and evolution and thereby provide insights into the role of phages in *S. enterica* taxonomy, diversity and biology.

Results

142 *Salmonella* phage genome sequences and patterns of variation

Complete genome sequences of *S. enterica* prophages were searched and downloaded from the NCBI database. Full genome sequences were available for 142 phages (Document S1) and their corresponding genomic information are summarized in Table 1 and include accession number, phage name, assigned cluster, host species, genome size, guanine plus cytosine (GC) content, number of ORFs and virus lineage and DNA structure, i.e., double stranded (dsDNA) or single stranded (ssDNA). The size range of the phage genomes was from 6.4-kb to 358.7-kb, with the majority between 30-kb to 50-kb (Fig. 1A), the GC content ranged from 35.5% to 65.4% (Table 1). The virus lineages for all 142 phages were retrieved from the Virus-Host DB (<https://www.genome.jp/virushostdb/>) and summarized in Table 1. Ninety-five percent of the phage genomes (135 out of 142) were linear ds DNA and belong to the order Caudovirales and four out of its five known families, namely: *Myoviridae*, *Siphoviridae*, *Podoviridae* and *Ackermannviridae* based on virus lineages retrieved from Virus-Host DB. There is a total of 28 genera represented in this collection of 142 prophages. Four of the remaining seven phages (5%) were single stranded DNA (NC_001954.1, NC_006294.1, NC_001332.1 and NC_025824.1), while three have not yet been classified (NC_010393.1, NC_010392.1 and NC_010391.1).

Open reading frame characterization of phage genomes

The availability of the 142 phage sequences in the NCBI database facilitated comparative genomic analysis. However, 32 out of 142 phages downloaded from the GenBank contained invalid start or stop codons for some ORFs, which were detected during our construction of the *Salmonella* prophage database (SpDB) and analysis with the Phamerator software (see under Materials and Methods). To ensure congruence between the annotations shown in the GenBank and ORFs displayed by the Phamerator, it became necessary to ensure that proper start and stop codons were present in the sequences. The detailed error messages (including number of errors and their locations in the original sequences) are shown in Table 1, and the revised sequences and NCBI annotation files are now included in Document S2. The distribution of the genome sizes mirrored the number of ORFs, with the genome size (grey) matching the number of ORFs (blue) as displayed in Fig. 1A and 1B. For instance, the 4 genomes with the smallest size (6408, 6744, 7107 and 8454 bp) had the least ORFs (10, 9, 12, and 10, respectively). Similarly, the 10 largest genomes encoded the highest number of ORFs, typically over 120 ORFs (Table 1A and 1B). There was a statistically significant, strong linear correlation between the genome sizes and number of ORFs ($R^2 = 0.95$, $p < 0.001$, Fig. 1C).

Salmonella phages occur in other bacteria

Although the 142 prophages were identified in *Salmonella enterica* strains present in the Salfos database [17], many prophages matched sequences of viral origin associated with bacterial hosts other than *Salmonella*. This designation of a non-*Salmonella* host was presumably a consequence of which host the prophage was associated with at the time of initial documentation or publication. The original known host lineage for each phage was retrieved online from Virus-Host Database (<https://www.genome.jp/virushostdb/>), which was used to evaluate the occurrence of these phages in other bacteria. As shown in Table 1 and illustrated in Fig. 2, fifty-three out of the 142 *Salmonella* phages (37.3%) were apparently first recovered from the genus *Escherichia*, followed by 34 phages (23.9%) first described for a *Salmonella* host. The others, including *Shigella*, *Burkholderia*, and *Pseudomonas*, showed relatively lower frequencies of 9, 6, and 6 phages, respectively (Fig. 2). Although the cellular host for the phage P4 is named as *Escherichia*, it is indeed a satellite virus for another phage called *Escherichia* virus P2, the latter serving as a helper to provide late gene functions for phage P4 lytic growth cycle, but not for its early functions especially DNA synthesis and lysogenization [37, 38]. We evaluated the above observations by using a web based tool called Hostphinder [36; <https://cge.cbs.dtu.dk/services/HostPhinder/>] and found a 97% agreement with the metadata on the bacterial host documented in the Virus-Host Database (Table S1).

Similarities among the 142 phage genomes based on nucleotide identity

Given that nucleotide identity and genome alignment are key tools for comparative genomic analysis and cluster assignment, NUCmer/MUMmer software was initially applied to analyze these 142 prophage sequences. The pairwise nucleotide identity was calculated among all the 142 genomes and those fragments with over 80% identity between two genomes were listed in Table S2. The sizes of aligned phage genome fragments varied, ranging from 103 bp to

14,505 bp. Out of the 142 genomes investigated, 133 share at least one fragment with another prophage. To illustrate the nucleotide connections between all the analyzed phage genomes, the visualization tool Circos [39] was used. *Salmonella*_phage_SJ46 (103 kb) and *Enterobacteria*_phage_P1 (95 kb), shared a large number of fragments with other *Salmonella* prophages as shown in Figure 3. In a striking contrast, *Salmonella/Cronobacter* prophage vB_CsaM_GAP32 and *Salmonella/cyanophage* MED4–213, which have the two biggest genomes (181- and 359-kb) did not share any fragment with another phage genome.

Clustering of phage genomes

Conserved DNA fragments among groups of prophage sequences (Fig. 3), were combined with both the results with ANI, identified with the aid of Kalign [31] and whole genome dot plot analysis, to assign the prophage genomes to clusters. To this end, a phylogenetic tree was constructed using MEGA X from the genome nucleotide identity matrix generated with the Kalign algorithm (Figure. S1). Furthermore, all 142 genomes were concatenated into a single nucleotide sequence and duplicated to form two axes for the purpose of generating a dot plot matrix (Fig. 4). We were able to assign 90 phage genomes into 17 clusters, named A to Q as follows: Cluster A (n = 3), Cluster B (n = 5), Cluster C (n = 2), Cluster D (n = 15), Cluster E (n = 4), Cluster F (n = 9), Cluster G (n = 5), Cluster H (n = 10), Cluster I (n = 4), Cluster J (n = 6), Cluster K (n = 12), Cluster L (n = 3), Cluster M (n = 3), Cluster N (n = 3), Cluster O (n = 2), Cluster P (n = 2) and Cluster Q (n = 2). The remaining 52 phage genomes could not be assigned to any cluster and remained as singletons. We observed both qualitative and quantitative differences in the structure of the clusters based on the intensity of the dot plots (Fig. 4) and pairwise nucleotide similarity between members of each cluster (Table 2, Cluster A-Q). Clusters E, F, H, I and J had relatively high intracluster nucleotide similarities and moderate genome sizes (37 - 77 kb). All four members of Cluster E belonged to the same genus, Epsilon15 virus under the family of *Podoviridae* according to the International Committee on Taxonomy of Viruses (ICTV) classification. Details of cluster assignment for all prophages are shown in Table 1.

We observed uniformity among the genome sizes and number of ORFs of members of the same cluster (Fig. 1C) which underscores the nucleotide identity among related genomes as also shown in Fig. 3.

Genome maps of multiple phages that incorporate and display nucleotide and amino acid sequence relationships

Using a ClustalW threshold of 35% amino acid identity and a BLASTP score of $1e-50$, the predicted ORFs and translated nucleotide sequences were assigned to groups of closely related sequences using the Phamerator software (Document S3 and Fig. 5). A total of 5796 Phamilies was assigned by Phamerator (Table S3). The most common Phamily was present in 49 prophages but there were 4330 Phamilies found in only one prophage. The relatively conserved Phamily numbers were summarized in the 17 assigned clusters in Table 3. To establish cluster-specific markers, we retrieved the conserved phamilies from each analyzed clusters and found that a total of 181 representative protein groups were present in all 17 clusters and 159 of them (excluding the 22 red highlighted proteins in Table 3) were specifically present in one cluster. For example, Cluster A uniquely contained seven Phamilies. In contrast, Cluster H contained 10 Phamilies but not all were unique because two of these Phamilies were also present in Cluster I. In the same vein, Cluster K contained 15 Phamilies, seven of which were shared with Cluster L. Thus, we demonstrated the presence of unique proteins and/or unique combination of proteins that define each prophage cluster, notwithstanding the fact that some individuals' proteins may be shared among some clusters. A representative genomic map of phages in Cluster H is shown in Fig. 5. Considerable genome length was observed to be conserved between members of the same cluster inferring synteny (violet shading blocks), with the same phamily ORF (same colour, Fig. 5). Often syntenic regions are interspersed with dissimilar and variable sequences (white blocks or breaks).

Discussion

We have carried out a comparative genomic analysis for the purpose of characterizing the prophages of *Salmonella enterica*. Both dsDNA and ssDNA viruses were represented in our collection of 142 phage genomes. The four ssDNA phages present in our collection belonged to the family *Inoviridae*. In contrast, the dsDNA phages were spread over four of the five known families of the order Caudovirales, i.e., *Myoviridae*, *Podoviridae*, *Siphoviridae* and the rare *Ackermannviridae*. Within these four families, a total of 28 different phage genera were represented (Table 1). Earlier studies using core genes analysis indicated that *Salmonella* phages could be classified into five groups, namely: P27-like, P2-like, lambdoid, P22-like, and T7-like [9, 10], and all of which were present in our prophage collection. From our classification, we have identified two new members of Cluster D namely, ST64T and ST104 which are related to the previously described P22-like group. We have described an additional 13 members in this group (Table 1). Similarly, we detected the P2-like PSP3 phages and were able to cluster them with an additional 12 double stranded phage viruses to make up Cluster K. In addition, three lambdoid phages, namely Gifsy 1, Gifsy 2 and lambda were assigned to lambdoid phage group Cluster M (Table 1 and Table 3). This work has extended published observations by identifying additional members of previously described, albeit small groupings, and has achieved a more discriminative and extensive characterization of *Salmonella* prophage sequences.

An earlier genomic comparison of tailed phages showed 337 fully sequenced lytic and temperate phages in the entire Enterobacteriaceae family [28], and based on this observation, a large number of diverse phages could potentially infect *Salmonella*. We observed the presence of the same phages infecting different bacteria and whether this is an outcome of the shared location or relatedness among hosts cannot be ascertained at this time. It is possible that both phylogeny, i.e., the relatedness among hosts such as belonging to the same family, or occupation of the same niche, i.e., gastrointestinal tract location may facilitate the presence of same prophages in different hosts. As examples, we observed phages X29 and KSF–1phi in *Salmonella*, which were first found in *Vibrio cholerae* according to Virus-Host DB [TAX:666; <https://www.genome.jp/virushostdb/>]. On the other hand, 38 other phages known to infect *Vibrio cholera* have not been reportedly found in *S. enterica* and given that the two organisms belong to different Orders, this suggests that hosts phylogeny rather than co-location plays the primary role whether prophages are shared among hosts. Nevertheless, it is difficult to entirely discount the role of a shared niche since the virus will still have to find the new host before infection can take place. Furthermore, 33 phages analyzed here were observed to have originated from

Escherichia coli strains [TAX:562] (Table 1). *Enterobacteria* phage fiAA91-ss is also able to infect at least two more hosts, namely, *Shigella sonnei* [TAX:624] and *Escherichia coli* O157:H7 [TAX:83334]. *Haemophilus* phage Aaphi23 can also infect *Aggregatibacter actinomycetemcomitans* [TAX:714] and *Haemophilus* [TAX:724]. The species *A. actinomycetemcomitans* has now been renamed *Haemophilus actinomycetemcomitans* by Potts *et al.* (1985) [40]. Based on our observations, studies of phage host range should not be restricted to specific species but should comprehensively involve as many different host genera as possible to capture all available information, even if the focus is a particular host species. This will help provide a broader perspective of the distribution of phages and contribute to their role in the evolution of the host.

The occurrence of the same phage sequences in different hosts may also imply horizontal viral gene transfer among hosts belonging to different genera. Genome clustering facilitates the identification of genes that are in greatest genetic flux and are more likely to have been exchanged horizontally during a relatively recent evolutionary time. Such viral sequence exchanges may help a phage increase its fitness to invade a new host, and evade selective pressure such as anti-phage defense mechanisms [11]. Given the biological arms race between bacteria and phages, and in order to thrive in most environments, phages have evolved multiple tactics to avoid, circumvent or subvert bacterial anti-phage mechanisms [21]. Ironically, these viral sequences once established in *Salmonella* may help the host to thrive in specific ecological niches, including the gut [41].

Diverse phage genomes were identified in our *Salmonella* phage collection. As shown in Fig. 2, the highest number of matching prophages were named after the genus *Escherichia* ($n = 53$) while *Salmonella* ranked second ($n = 34$). Regarding the lineage for their original known host, three phyla (Firmicutes, Proteobacteria and Cyanobacteria), four classes (Bacilli, Betaproteobacteria, Alphaproteobacteria and Gammaproteobacteria) and 25 unique genera could be identified (Table 1). Such a wide host span provides further evidence of the diversity of *Salmonella* prophages analyzed in this study. In a study of prophages integrated in a single host species *Mycobacterium smegmatis*, a threshold of 50% nucleotide identity was used for genome cluster assignment [24]. The threshold was slightly reduced (45%) for clustering *Pseudomonas* phages because phages infecting a genus would be expected to show greater variation in genome sequences than one infecting a single species [23]. Among the 56 phage clusters reported for the Enterobacteriaceae family, the sequence similarity was substantially less between clusters [28], indicating a higher degree of variation and justifying a lower threshold of nucleotide identity for certain clusters in *Salmonella* phages, a large proportion of which may infect or have previously infected other hosts.

It should be noted that nucleotide identity is not the only parameter for assessing genome properties, because the nucleotide alignments for thousands of homologous protein are not significant based on nucleotide alignment, but are clearly homologous based on statistically significant protein structural similarity or strong sequence similarity to an intermediate sequence [42]. Thus, there may not be a linear relationship between sequence identity and function [43]. In our set of phage genomes, except for Cluster B, L and M showed a lower pairwise ANI of 41%, all the other clusters Clusters E (59%), F (75%) and J (57%) displayed high nucleotide identity (Table 3). Their assignment to each of these clusters was supported by results of analysis using dotplot program, Kalign genome alignment and gene content analysis. For instance, the dotplot (Fig. 4) and Kalign analysis grouped members of Clusters B, G and N, even though some of their respective nucleotide identities were 40.7%, 42.2%, and 42.3% (Fig. 4 and Figure S1). A similar phenomenon was also observed for Cluster L made up of members belonging to the same P2 virus group showing a nucleotide identity of 41.3%. The differences in the output of the different tools should not be surprising because of their unique underlying algorithms. While Kalign focuses more on analyzing larger genomes in general, MUMmer focuses more on the similar DNA fragment identification. Despite the high degree of diversity in our prophage collection, we were still able to cluster related isolates using congruent results from at least two bioinformatics analyses.

The genome size ranges of the prophages documented for the different bacteria genera are fairly similar: *Salmonella* (6.4 - 358.7 kb), *Pseudomonas* (3.0 - 316.0 kb), *Staphylococcus* (15.6 - 138.7 kb), *Gordonia* (17.1 - 103.4 kb), *Bacillus* (14.3 - 497.5 kb) and *Mycobacterium* (41.9 - 164.6 kb). The ranges of the GC content showed less of an overlap: *Salmonella* (35.5 to 65.4%), *Pseudomonas* (37.0 to 66.0%), *Staphylococcus* (29.3 to 38.0%), *Gordonia* (47.0 to 68.8%), *Bacillus* (29.9 to 49.9%) and *Mycobacterium* (56.3 to 69.1%) [23–27]. *Salmonella* and *Pseudomonas* both belong to the Enterobacteriaceae family and their phages share very similar genome sizes and GC content. Despite the similarities between the phages of *Pseudomonas* and *Salmonella*, the former appear to display better clustering pattern (fewer singletons) based on the grouping of 100 out of 130 phages [23] compared to 90 out of 142 *Salmonella* phages with 52 singletons. However, as *Pseudomonas* bacteriophages were collected only using "*Pseudomonas*" as host for the search in the database [23], the set most likely did not represent the full complement of viruses capable of infecting *Pseudomonas* and integrating into the genome and would have excluded bacteriophages of this group but first found or described in another bacterial host. We expect that more diverse prophage patterns would be obtained for *Pseudomonas* and other bacterial hosts if a more comprehensive search of bacterial genomes is carried with tool such as PHASTER [34].

The diversity of *Salmonella* prophage genomes was also reflected in the total number of phamilies for the ORFs in the analyzed prophage genomes: 5796. One phamily with Pham number of 2217 was observed to be dominant and was present in 49 prophages (34.5% of 142 phages) whereas 4330 phamilies were each present in a single prophage, which makes it challenging to select some conserved genes for all the 142 prophage genomes. Clustering of the viral genome was useful in establishing relatedness of *Salmonella* bacteriophages. In each assigned cluster, some conserved Pham numbers (containing different ORFs) are present. For example, Pham 180 (portal protein), Pham 2012 (recombination protein) and Pham 2217 (endopeptidase) are commonly present in Cluster D; Pham 321 (phage head-tail connector protein), Pham 415 (terminase large subunit) and Pham 1522 (terminase small unit) in Cluster E; Pham 1995 (lysozyme), Pham 2370 (terminator) and Pham 1332 (attachment invasion locus protein precursor) in Cluster F; Pham 27 (phage tail protein), Pham 519 (phage portal protein), and Pham 1717 (assembly protein) in Cluster H; Pham 528 (major capsid protein), Pham 297 (terminase large subunit) and Pham 666 (tail protein) in Cluster J; Pham 963 (base plate assembly protein) in Cluster K (Document S3). Specifically, some proteins are unique to one cluster, for example, four members of Pham 4878 (a hypothetical protein), Pham 1893 (a hypothetical protein) Pham 2968 (a hypothetical protein) Pham 2849 (a hypothetical protein) in the Cluster E. These may be good markers for characterizing prophage members of the different clusters (Document S3, Table 3).

The observations reported in this study are quite relevant for the application of bacteriophages as antibacterial agents and in cloning vector construction. Our list of *Salmonella* bacteriophages can be used for screening a novel, candidate bacteriophage identified as a potential anti-bacterial agent for *Salmonella* or any host described in this study. The implication is that because the bacteriophages present in our collection induce lysogeny, the bacterial host will be

immune to infection or lysis by the same bacteriophage; a bacteriophage on our list will likely not be an effective antibacterial agent for the hosts identified in this study. Thus, a distinct bacteriophage may be a better anti-bacterial candidate than one on our list. Similarly, the *Salmonella* prophage database in the Phamerator can be used to evaluate a candidate antibacterial agent even if it is distinct from members on our list. Because bacteriophages are prone to recombination leading to a mosaic profile, the protein components can be used to assess relatedness with the goal of choosing a candidate antibacterial agent that is phylogenetically distant from any of the isolates in our collection to increase the chance of success. In the same vein, knowledge from our collection can be used in strategies to design phage vectors. For example, λ cloning vectors requires a lytic cycle and their ability to package large foreign DNA fragments have relied on the removal of lysogenic genes from the vectors. Thus, the removal of lysogenic fragments in a temperate phage can probably deviate the life cycle into a lytic path making them more relevant for vector construction especially if the bacteriophage has signature genetic markers that can be exploited for selection or vector purification, e.g., antibacterial resistance genes or a target for a widely used ligand.

Conclusions

The comparative genomic analysis of 142 *Salmonella enterica* subsp. *enterica* prophages revealed a high diversity in genomic characteristics, compared to that in other bacteria species such as *Pseudomonas*, *Staphylococcus*, *Gordonia*, *Bacillus* and *Mycobacterium*. The combination of nucleotide identity, dot plot, genome map comparison and gene content analysis, revealed the presence of 17 main clusters of *Salmonella* phages and many singletons. In order to have a fuller picture of *Salmonella* phages, a similar comparative phage genomic analysis needs to be performed on *Salmonella* virulent/lytic phages. The high diversity among prophages may well be a mechanism developed to generate new molecules and decoys to thwart the potent, anti-viral defence mechanism of the bacterial host. We hypothesize that in place of the resources needed to lyse a host cell, temperate prophages may instead have developed a rather sophisticated capacity to acquire and display diversity and thereby present a degree of invincibility against the host arsenal so that they can survive long enough to integrate into the host genome. Thus, we predict that prophages will show more diversity than their virulent phage counterparts. Areas of conservation and variations among the investigated prophage genomes provides further evidence showing why prophage typing is a discriminative method for *Salmonella* typing. A fuller understanding of the genomic architecture of *Salmonella* bacteriophages should furnish practical information relevant for bacterial strain construction, vector development, and the selection of appropriate phages to be tested for bio-control strategies.

Methods

Phage genome sequences

We previously identified 154 different prophages among 1,760 *S. enterica* genomes derived from 151 serovars using PHASTER [18]. We downloaded 142 of these 154 phage genomes from NCBI Batch Entrez (<https://www.ncbi.nlm.nih.gov/sites/batchentrez>) but were unable to locate the full length of the remaining 12 genomes. Genome annotation were downloaded from NCBI and validated using gene calling programs GeneMarkS and Glimmer [44–46], and BLASTN when necessary.

Comparative phage genome analysis

All 142 phage genome sequences were pooled and saved as a multi-fasta file and aligned to one another using MUMmer v4.0.0.beta2. Genome comparison was carried out to produce delta files using the following parameters: maxgap = 200; mincluster = 90; minmatch = 60. Results were generated as coordinate files using “shwon-coords” and visualized via Circos [39]. Whole genome alignment and calculation of percentage of nucleotide identity were carried out with Kalign [30]. The evolutionary history was inferred using the Neighbor-Joining method [47]. The bootstrap consensus tree inferred from 500 replicates [48] was taken to represent the evolutionary history of the taxa analyzed. Branches corresponding to partitions reproduced in less than 50% bootstrap replicates were collapsed. The percentage of replicate trees in which the associated taxa clustered together in the bootstrap test (500 replicates) was shown next to the branches. The evolutionary distances were computed using the Maximum Composite Likelihood method [49] and are presented as the number of base substitutions per site. There were a total of 431,295 positions in the final dataset. Evolutionary analyses for tree construction were conducted in MEGA X [50]. Prophage genomes (n = 142 phage) were concatenated into a single sequence with a total length of 7,260,982 bp, which when plotted against itself with a sliding window of 10 bp and visualized by Gepard 1.40 version [29], revealed an overall pattern of similarity or dissimilarity of all the genomes. The graphics displayed pairwise similarity between genomes which was then used for the preliminary assignment of clusters. Among all the analyzed prophage genomes, if two sequences shared high similarity, a diagonal would show at that location on the plot (the center diagonal line demonstrated the 100% similarities where a sequence was compared to itself).

Genome clustering

Three criteria were used to cluster the phage genomes. First, the genomes were grouped based on nucleotide identity among members. Second, dot plot was used to analyze sequences based on similarity leading to graphically demonstrable clustering of sequences. Third, translated nucleotide sequences were used to cluster phages based on translated amino acid sequences. Phage genomes that did not meet these criteria were identified as ‘singletons’.

Salmonella phage database creation and genome map viewing via Phamerator

In order to produce the first, web-based inventory of *Salmonella* prophages that could be used for comparative analysis with prophage genomes from other bacteria, we created the SpDB in the Phamerator platform. For this purpose, *Salmonella* phage database, a web-based application PhamDB was used for building the *Salmonella* Phamerator phage database consisting of 142 phage sequences. Briefly, after installing Docker Toolbox, Kitematic was launched to finish the initial setup and loading. An existing ‘PhamDB’ database in the Phamerator platform was downloaded and used as a template. By running the

PhamDB program as a web interface on a local network, a new database was created in toolbar using GenBank Files as input. All the 142 phage NCBI files were summarized in Document S2. The generated database was a sql file which was used as an input file and uploaded into Phamerator website (<https://phamerator.org>, created and maintained by Dr. Steven Cresawn of James Madison University). Based on the assigned clusters, genome maps can be visualized for direct comparisons. As displayed in the Phamerator map, long regions of violet shading indicate long conserved regions between phage genomes. Within a cluster, the same color block represents the ORF with higher similarities. Regions of high similarity and same-coloured ORF blocks shown on the map indicated a prevalent synteny. Areas with little or no sequence similarity between genome sequences are shown as either white blocks or a break in a syntenic block.

List Of Abbreviations

ORFs: Open Reading Frames

PHASTER: PHAge Search Tool Enhanced Release

CRISPRs: clustered regularly interspaced short palindromic repeats

ICTV: International Committee on Taxonomy of Viruses

Gepard: Genome Pair Rapid Dotter

Double stranded DNA (dsDNA)

Single stranded DNA (ssDNA)

Declarations

Ethics approval and consent to participate

Not applicable

Consent for publication

This manuscript was approved for publication by the Canadian Food Inspection Agency?]

Availability of data and material

Not applicable

Competing interests

The authors declare that they have no competing interests

Funding

RG was supported by the Genome Canada funded project titled "A Syst-OMICS approach to ensuring food safety and reducing the economic burden of Salmonellosis". DO's research program has received funding support from Genome Research and Development Initiative of the Government of Canada, Ontario Ministry of Agriculture, Food and Rural Affairs, Canadian Security and Science Program of the Department of National Defence and the Canadian Food Inspection Agency. SM holds a Tier 1 Canada Research Chair in Bacteriophages. LG is the Leung Family Professor of Food Safety in University of Guelph, Guelph.

Authors' contributions

RG and DO conceived and designed the study. RG developed and analyzed the data and wrote the draft manuscript. SN, SM analyzed data and edited manuscript. LG secured funding and edited manuscript, DO secured funding, supervised the project, analyzed data and edited manuscript. All authors read and approved the final manuscript.

Availability of data and materials

The datasets supporting the conclusions of this article are included within the article and its additional files.

Acknowledgements

We thank Dr. Walid Mottawea of the University of Ottawa, Canada for assistance in prophage sequence analysis, and Dr. Steven Cresawn from James Madison University for assistance with uploading the *Salmonella* phage Database on the Phamerator website (<https://phamerator.org>).

We acknowledge excellent review comments by Ray Theoret, Hongsheng Huang and Amit Matthews.

References

1. Fabrega A, Vila J: *Salmonella enterica* serovar Typhimurium skills to succeed in the host: virulence and regulation. *Clin Microbiol Rev* 2013, 26(2):308-341.
2. Su LH, Chiu CH: *Salmonella*: clinical importance and evolution of nomenclature. *Chang Gung Med J* 2007, 30(3):210-219.
3. Gal-Mor O, Boyle EC, Grassl GA: Same species, different diseases: how and why typhoidal and non-typhoidal *Salmonella enterica* serovars differ. *Front Microbiol* 2014, 5:391.
4. Jantsch J, Chikkaballi D, Hensel M: Cellular aspects of immunity to intracellular *Salmonella enterica*. *Immunol Rev* 2011, 240(1):185-195.
5. Hohmann EL: Nontyphoidal salmonellosis. *Clin Infect Dis* 2001, 32(2):263-269.
6. Ryan KJaR, C.G., Eds.: *Sherris medical microbiology*. 4th Edition, McGraw-Hill, New York 2004.
7. Lin DM, Koskella B, Lin HC: Phage therapy: An alternative to antibiotics in the age of multi-drug resistance. *World J Gastrointest Pharmacol Ther* 2017, 8(3):162-173.
8. Wommack KE, Colwell RR: Virioplankton: viruses in aquatic ecosystems. *Microbiol Mol Biol Rev* 2000, 64(1):69-114.
9. Kropinski AM, Sulakvelidze A, Konczy P, Poppe C: *Salmonella* phages and prophages—genomics and practical aspects. *Methods Mol Biol* 2007, 394:133-175.
10. Switt AI, Sulakvelidze A, Wiedmann M, Kropinski AM, Wishart DS, Poppe C et al: *Salmonella* phages and prophages: genomics, taxonomy, and applied aspects. *Methods Mol Biol* 2015, 1225:237-287.
11. Labrie SJ, Samson JE, Moineau S: Bacteriophage resistance mechanisms. *Nat Rev Microbiol* 2010, 8(5):317-327.
12. Doron S, Melamed S, Ofir G, Leavitt A, Lopatina A, Keren M et al: Systematic discovery of antiphage defense systems in the microbial pangenome. *Science* 2018, 359(6379).
13. Levesque S, Moineau S: A stockpile of antiviral defences. *Nature* 2018, 556(7701):318-319.
14. Ackermann HW: Tailed bacteriophages: the order caudovirales. *Adv Virus Res* 1998, 51:135-201.
15. Jakub Barylski FE, Bas E, Dutilh, Margo B.P. Schuller, Robert A. Edwards, Annika Gillis, Jochen Klumpp, Petar Knezevic, Mart Krupovic, Jens H. Kuhn, Rob Lavigne, Hanna M. Oksanen, Matthew B. Sullivan, Johannes Wittmann, Igor Tolstoy, J. Rodney Brister, Andrew M. Kropinski, Evelien M Adriaenssens: Genomic, proteomic, and phylogenetic analysis of spounaviruses indicates paraphyly of the order Caudovirales *BioRxiv* 2017.
16. Duckworth DH: "Who discovered bacteriophage?". *Bacteriol Rev* 1976, 40(4):793-802.
17. Boyd JS: The symbiotic bacteriophages of *Salmonella typhimurium*. *J Pathol Bacteriol* 1950, 62(4):501-517.
18. Mottawea W, Duceppe MO, Dupras AA, Usongo V, Jeukens J, Freschi L et al: *Salmonella enterica* prophage sequence profiles reflect genome diversity and can be used for high discrimination subtyping. *Front Microbiol* 2018, 9:836.
19. Brussow H, Canchaya C, Hardt WD: Phages and the evolution of bacterial pathogens: from genomic rearrangements to lysogenic conversion. *Microbiol Mol Biol Rev* 2004, 68(3):560-602, table of contents.
20. Colavecchio A, D'Souza Y, Tompkins E, Jeukens J, Freschi L, Emond-Rheault JG et al: Prophage integrase typing is a useful indicator of genomic diversity in *Salmonella enterica*. *Front Microbiol* 2017, 8:1283.
21. Nicolle P, Vieu JF, Diverneau G: Supplementary lysotyping of Vi-positive strains of *Salmonella typhi*, insensitive to all the adapted preparations of Craigie's Vi II phage (group I+IV). *Arch Roum Pathol Exp Microbiol* 1970, 29(4):609-617.
22. Anderson ES, Ward LR, Saxe MJ, de Sa JD: Bacteriophage-typing designations of *Salmonella typhimurium*. *J Hyg (Lond)* 1977, 78(2):297-300.
23. Ha AD, Denver DR: Comparative genomic analysis of 130 bacteriophages infecting bacteria in the genus *Pseudomonas*. *Front Microbiol* 2018, 9:1456.
24. Hatfull GF, Jacobs-Sera D, Lawrence JG, Pope WH, Russell DA, Ko CC et al: Comparative genomic analysis of 60 Mycobacteriophage genomes: genome clustering, gene acquisition, and gene size. *J Mol Biol* 2010, 397(1):119-143.
25. Kwan T, Liu J, DuBow M, Gros P, Pelletier J: The complete genomes and proteomes of 27 *Staphylococcus aureus* bacteriophages. *Proc Natl Acad Sci U S A* 2005, 102(14):5174-5179.
26. Grose JH, Jensen GL, Burnett SH, Breakwell DP: Genomic comparison of 93 *Bacillus* phages reveals 12 clusters, 14 singletons and remarkable diversity. *BMC Genomics* 2014, 15:855.
27. Pope WH, Mavrich TN, Garlena RA, Guerrero-Bustamante CA, Jacobs-Sera D, Montgomery MT et al: Bacteriophages of *Gordonia* spp. Display a spectrum of diversity and genetic relationships. *MBio* 2017, 8(4).

28. Grose JH, Casjens SR: Understanding the enormous diversity of bacteriophages: the tailed phages that infect the bacterial family Enterobacteriaceae. *Virology* 2014, 468-470:421-443.
29. Krumsiek J, Arnold R, Rattei T: Gepard: a rapid and sensitive tool for creating dotplots on genome scale. *Bioinformatics* 2007, 23(8):1026-1028.
30. Lassmann T, Sonnhammer EL: Kalign—an accurate and fast multiple sequence alignment algorithm. *BMC Bioinformatics* 2005, 6:298.
31. Kurtz S, Phillippy A, Delcher AL, Smoot M, Shumway M, Antonescu C et al: Versatile and open software for comparing large genomes. *Genome Biol* 2004, 5(2):R12.
32. Cresawn SG, Bogel M, Day N, Jacobs-Sera D, Hendrix RW, Hatfull GF: Phamerator: a bioinformatic tool for comparative bacteriophage genomics. *BMC Bioinformatics* 2011, 12:395.
33. Lamine JG, DeJong RJ, Nelesen SM: PhamDB: a web-based application for building Phamerator databases. *Bioinformatics* 2016, 32(13):2026-2028.
34. Arndt D, Grant JR, Marcu A, Sajed T, Pon A, Liang Y et al: PHASTER: a better, faster version of the PHAST phage search tool. *Nucleic Acids Res* 2016, 44(W1):W16-21.
35. Zhou Y, Liang Y, Lynch KH, Dennis JJ, Wishart DS: PHAST: a fast phage search tool. *Nucleic Acids Res* 2011, 39(Web Server issue):W347-352.
36. Villarroel J, Kleinheinz KA, Jurtz VI, Zschach H, Lund O, Nielsen M et al: HostPhinder: A Phage Host Prediction Tool. *Viruses* 2016, 8(5).
37. Haggard-Ljungquist E, Jacobsen E, Rishovd S, Six EW, Nilssen O, Sunshine MG et al: Bacteriophage P2: genes involved in baseplate assembly. *Virology* 1995, 213(1):109-121.
38. Six EW: The helper dependence of satellite bacteriophage P4: which gene functions of bacteriophage P2 are needed by P4? *Virology* 1975, 67(1):249-263.
39. Krzywinski M, Schein J, Birol I, Connors J, Gascoyne R, Horsman D et al: Circos: an information aesthetic for comparative genomics. *Genome Res* 2009, 19(9):1639-1645.
40. Potts TV, Zambon, J. J., and Genco R. J.: Reassignment of *Actinobacillus actinomycetemcomitans* to the Genus *Haemophilus* as *Haemophilus actinomycetemcomitans* comb. nov. *International Journal of Systematic Bacteriology* 1985, 35(3):337-341.
41. Shkoporov AN, Hill C: Bacteriophages of the human gut: The "known unknown" of the microbiome. *Cell Host Microbe* 2019, 25(2):195-209.
42. Pearson WR: An introduction to sequence similarity ("homology") searching. *Curr Protoc Bioinformatics* 2013, Chapter 3:Unit3 1.
43. Joshi T, Xu D: Quantitative assessment of relationship between sequence similarity and function similarity. *BMC Genomics* 2007, 8:222.
44. Besemer J, Borodovsky M: Heuristic approach to deriving models for gene finding. *Nucleic Acids Res* 1999, 27(19):3911-3920.
45. Delcher AL, Harmon D, Kasif S, White O, Salzberg SL: Improved microbial gene identification with GLIMMER. *Nucleic Acids Res* 1999, 27(23):4636-4641.
46. Mills R, Rozanov M, Lomsadze A, Tatusova T, Borodovsky M: Improving gene annotation of complete viral genomes. *Nucleic Acids Res* 2003, 31(23):7041-7055.
47. Saitou N, Nei M: The neighbor-joining method: a new method for reconstructing phylogenetic trees. *Mol Biol Evol* 1987, 4(4):406-425.
48. Felsenstein J: Confidence limits on phylogenies: an approach using the bootstrap. *Evolution* 1985, 39(4):783-791.
49. Tamura K, Nei M, Kumar S: Prospects for inferring very large phylogenies by using the neighbor-joining method. *Proc Natl Acad Sci U S A* 2004, 101(30):11030-11035.
50. Kumar S, Stecher G, Li M, Knyaz C, Tamura K: MEGA X: molecular evolutionary genetics analysis across computing platforms. *Mol Biol Evol* 2018, 35(6):1547-1549.

Tables

Table 1 The profiles of 142 prophages present in *Salmonella enterica*

Accession number	Phage name	Family	Genus	Cluster	Size (bp)	GC (%)	ORF	Lineage of original host	
								Family	Genus
NC_006552.1	<i>Pseudomonas</i> phage F116	<i>Podoviridae</i>	F116virus	A	65195	63.17	70	Pseudomonadaceae	<i>Pseudomonas</i>
NC_005357.1	<i>Bordetella</i> phage BPP-1	<i>Podoviridae</i>	Bpp1virus	A	42493	65.41	49	Alcaligenaceae	<i>Bordetella</i>
NC_005887.1	<i>Burkholderia</i> phage BcepC6B	<i>Podoviridae</i>	Bpp1virus	A	42415	65.19	46	Burkholderiaceae	<i>Burkholderia</i>
NC_015266.1	<i>Burkholderia</i> phage KL3	<i>Podoviridae</i>	P2virus	B	40555	63.23	52	Burkholderiaceae	<i>Burkholderia</i>
NC_025115.1	<i>Ralstonia</i> phage RSY1 DNA	<i>Myoviridae</i>	P2virus	B	40002	64.82	49	Burkholderiaceae	<i>Ralstonia</i>
NC_015273.1	<i>Burkholderia</i> phage KS14	<i>Myoviridae</i>	P2virus	B	32317	62.28	44	Burkholderiaceae	<i>Burkholderia</i>
NC_009237.1	<i>Burkholderia</i> phage phiE255 chromosome	<i>Myoviridae</i>	Bcepmyovirus	B	37446	63.05	55	Burkholderiaceae	<i>Burkholderia</i>
NC_005882.1	<i>Burkholderia cenocepacia</i> phage BcepMu	<i>Myoviridae</i>	Bcepmyovirus	B	36748	62.86	53	Burkholderiaceae	<i>Burkholderia</i>
NC_005178.1	<i>Pseudomonas</i> phage D3112	<i>Siphoviridae</i>	D3112virus	C	37611	64.34	55	Pseudomonadaceae	<i>Pseudomonas</i>
NC_008717.1	<i>Pseudomonas</i> phage DMS3	<i>Siphoviridae</i>	D3112virus	C	36415	64.26	52	Pseudomonadaceae	<i>Pseudomonas</i>
NC_011976.1	<i>Salmonella</i> phage epsilon34	<i>Podoviridae</i>	P22virus	D	43016	47.26	73	Enterobacteriaceae	<i>Salmonella</i>
NC_030919.1	<i>Salmonella</i> phage 118970_sal4	<i>Podoviridae</i>	P22virus	D	42418	46.81	64	Enterobacteriaceae	<i>Salmonella</i>
NC_031019.1	<i>Enterobacteria</i> phage UAB_Phi20	<i>Podoviridae</i>	P22virus	D	41809	47.24	80	Enterobacteriaceae	<i>Salmonella</i>
NC_005841.1	<i>Enterobacteria</i> phage ST104 DNA	<i>Podoviridae</i>	P22virus	D	41391	47.43	63	Enterobacteriaceae	<i>Salmonella</i>
NC_028696.2	<i>Salmonella</i> phage SEN22	<i>Podoviridae</i>	P22virus	D	41338	47.83	55	Enterobacteriaceae	<i>Salmonella</i>
NC_014900.1	<i>Salmonella</i> phage ST160	<i>Podoviridae</i>	P22virus	D	40986	47.06	63	Enterobacteriaceae	<i>Salmonella</i>
NC_013059.1	<i>Salmonella</i> phage c341	<i>Podoviridae</i>	P22virus	D	40975	47.4	67	Enterobacteriaceae	<i>Salmonella</i>
NC_004348.1	<i>Enterobacteria</i> phage ST64T	<i>Podoviridae</i>	P22virus	D	40679	47.52	65	Enterobacteriaceae	<i>Salmonella</i>
NC_031946.1	<i>Salmonella</i> Phage 103203_sal5	<i>Podoviridae</i>	P22virus	D	40443	46.52	60	Enterobacteriaceae	<i>Salmonella</i>
NC_017985.1	<i>Salmonella</i> phage SPN9CC	<i>Podoviridae</i>	P22virus	D	40128	47.33	62	Enterobacteriaceae	<i>Salmonella</i>
NC_018275.1	<i>Salmonella</i> phage vB_SemP_Emek	<i>Podoviridae</i>	P22virus	D	39783	47.65	70	Enterobacteriaceae	<i>Salmonella</i>
NC_019501.1	<i>Enterobacteria</i> phage IME10	<i>Podoviridae</i>	P22virus	D	39646	47.5	53	Enterobacteriaceae	<i>Escherichia</i>
NC_005344.1	<i>Enterobacteria</i> phage Sf6	<i>Podoviridae</i>	P22virus	D	39043	47.47	66	Enterobacteriaceae	<i>Shigella</i>
NC_027398.1	<i>Enterobacteria</i> phage Sf101	<i>Podoviridae</i>	P22virus	D	38742	47.44	66	Enterobacteriaceae	<i>Shigella</i>
NC_002730.1	<i>Enterobacteria</i> phage HK620	<i>Podoviridae</i>	P22virus	D	38297	46.69	58	Enterobacteriaceae	<i>Escherichia</i>
NC_019445.1	<i>Escherichia</i> phage TL-2011b	<i>Podoviridae</i>	Epsilon15virus	E	44784	47.05	57	Enterobacteriaceae	<i>Escherichia</i>
NC_031077.1	<i>Enterobacter</i> phage Tyron	<i>Podoviridae</i>	Epsilon15virus	E	41760	50.59	56	Enterobacteriaceae	<i>Enterobacter</i>
NC_004775.2	<i>Enterobacteria</i> phage epsilon15	<i>Podoviridae</i>	Epsilon15virus	E	39672	50.83	51	Enterobacteriaceae	<i>Salmonella</i>
NC_016761.1	<i>Salmonella</i> phage SPN1S	<i>Podoviridae</i>	Epsilon15virus	E	38684	50.16	52	Enterobacteriaceae	<i>Salmonella</i>
NC_028656.1	<i>Enterobacteria</i> phage VT2phi_272	<i>Podoviridae</i>	Tl2011virus	F	65955	50.11	83	Enterobacteriaceae	<i>Escherichia</i>
NC_010237.1	<i>Enterobacteria</i> phage Min27	<i>Podoviridae</i>	Nona33virus	F	63395	49.5	83	Enterobacteriaceae	<i>Escherichia</i>
NC_028685.1	<i>Shigella</i> phage Ss-VASD	<i>Podoviridae</i>	Tl2011virus	F	62851	50.07	74	Enterobacteriaceae	<i>Shigella</i>
NC_025434.1	<i>Shigella</i> phage POCJ13	<i>Podoviridae</i>	Pocjvirus	F	62699	49.35	79	Enterobacteriaceae	<i>Shigella</i>
NC_000924.1	<i>Enterobacteria</i> phage 933W	<i>Podoviridae</i>	Nona33virus	F	61670	49.37	80	Enterobacteriaceae	<i>Escherichia</i>
NC_000902.1	<i>Enterobacteria</i> phage VT2-Sakai	<i>Podoviridae</i>	Nona33virus	F	60942	49.91	83	Enterobacteriaceae	<i>Escherichia</i>
NC_018846.1	<i>Escherichia</i> phage P13374	<i>Podoviridae</i>	Tl2011virus	F	60894	50.23	79	Enterobacteriaceae	<i>Escherichia</i>
NC_029120.1	<i>Shigella</i> phage 75_02 Stx	<i>Podoviridae</i>	Pocjvirus	F	60875	49.12	76	Enterobacteriaceae	<i>Shigella</i>
NC_008464.1	Stx2-converting phage 86	<i>Podoviridae</i>	Nona33virus	F	60238	49.07	81	Enterobacteriaceae	<i>Escherichia</i>
NC_004813.1	<i>Enterobacteria</i> phage BP-4795	<i>Siphoviridae</i>	unclassified	G	57930	50.61	85	Enterobacteriaceae	<i>Escherichia</i>
NC_011356.1	<i>Enterobacteria</i> phage YYZ-2008	<i>Siphoviridae</i>	unclassified	G	54896	51.12	75	Enterobacteriaceae	<i>Escherichia</i>
NC_011357.1	Stx2-converting phage 1717	<i>Siphoviridae</i>	unclassified	G	62147	50.92	77	Enterobacteriaceae	<i>Escherichia</i>
NC_018279.1	<i>Salmonella</i> phage vB_SosS_Oslo	<i>Siphoviridae</i>	unclassified	G	49116	48.74	79	Enterobacteriaceae	<i>Salmonella</i>
NC_006949.1	<i>Enterobacteria</i> phage ES18	<i>Siphoviridae</i>	unclassified	G	46900	48.59	79	Enterobacteriaceae	<i>Salmonella</i>
NC_019721.1	<i>Enterobacteria</i> phage mEp390	<i>Siphoviridae</i>	Hk97virus	H	40029	51.68	59	Enterobacteriaceae	<i>Escherichia</i>
NC_019705.1	<i>Enterobacteria</i> phage mEpX2	<i>Siphoviridae</i>	Hk97virus	H	38759	50.08	67	Enterobacteriaceae	<i>Escherichia</i>
NC_016160.1	<i>Escherichia</i> phage HK75	<i>Siphoviridae</i>	Hk97virus	H	36661	50.19	58	Enterobacteriaceae	<i>Escherichia</i>
NC_019709.1	<i>Enterobacteria</i> phage mEpX1	<i>Siphoviridae</i>	Hk97virus	H	41567	49.31	66	Enterobacteriaceae	<i>Escherichia</i>
NC_019719.1	<i>Enterobacteria</i> phage HK633	<i>Siphoviridae</i>	Hk97virus	H	41528	49.65	67	Enterobacteriaceae	<i>Escherichia</i>
NC_019714.1	<i>Enterobacteria</i> phage HK446	<i>Siphoviridae</i>	Hk97virus	H	39026	50.1	60	Enterobacteriaceae	<i>Escherichia</i>
NC_019708.1	<i>Enterobacteria</i> phage mEp235	<i>Siphoviridae</i>	Hk97virus	H	37595	50.01	61	Enterobacteriaceae	<i>Escherichia</i>
NC_002166.1	Bacteriophage HK022	<i>Siphoviridae</i>	Hk97virus	H	40751	49.48	65	Enterobacteriaceae	<i>Escherichia</i>
NC_002167.1	<i>Enterobacteria</i> phage HK97	<i>Siphoviridae</i>	Hk97virus	H	39732	49.79	61	Enterobacteriaceae	<i>Escherichia</i>
NC_019768.1	<i>Enterobacteria</i> phage HK106	<i>Siphoviridae</i>	Hk97virus	H	41468	49.34	65	Enterobacteriaceae	<i>Escherichia</i>
NC_021190.1	<i>Enterobacteria</i> phage phi80	<i>Siphoviridae</i>	unclassified	I	46150	52.13	63	Enterobacteriaceae	<i>Escherichia</i>
NC_019717.1	<i>Enterobacteria</i> phage HK225	<i>Siphoviridae</i>	unclassified	I	45366	51.96	69	Enterobacteriaceae	<i>Escherichia</i>
NC_019704.1	<i>Enterobacteria</i> phage mEp237	<i>Siphoviridae</i>	unclassified	I	44375	51.43	63	Enterobacteriaceae	<i>Escherichia</i>
NC_019706.1	<i>Enterobacteria</i> phage mEp043 c-1	<i>Siphoviridae</i>	unclassified	I	42780	50.79	69	Enterobacteriaceae	<i>Escherichia</i>
NC_031940.1	<i>Salmonella</i> phage 118970_sal3	<i>Myoviridae</i>	unclassified	J	77375	50.7	135	Enterobacteriaceae	<i>Salmonella</i>

NC_003356.1	<i>Enterobacteria</i> phage phiP27	<i>Myoviridae</i>	unclassified	J	42575	49.35	58	Enterobacteriaceae	Escherichia
NC_021857.1	<i>Shigella</i> phage SIII	<i>Myoviridae</i>	unclassified	J	41475	49.17	58	Enterobacteriaceae	Shigella
NC_004313.1	<i>Salmonella</i> phage ST64B	<i>Myoviridae</i>	unclassified	J	40149	51.01	56	Enterobacteriaceae	Salmonella
NC_022749.1	<i>Shigella</i> phage SIIIV	<i>Myoviridae</i>	unclassified	J	39758	50.3	54	Enterobacteriaceae	Shigella
NC_003444.1	<i>Enterobacteria</i> phage Sfv	<i>Myoviridae</i>	unclassified	J	37074	50.77	53	Enterobacteriaceae	Shigella
NC_001895.1	<i>Enterobacteria</i> phage P2	<i>Myoviridae</i>	P2virus	K	33593	50.17	43	Enterobacteriaceae	Escherichia
NC_004745.1	<i>Yersinia</i> phage L-413C	<i>Myoviridae</i>	P2virus	K	30728	52.11	40	Yersiniaceae	Yersinia
NC_005340.1	<i>Enterobacteria</i> phage PsP3	<i>Myoviridae</i>	P2virus	K	30636	52.83	42	Enterobacteriaceae	Salmonella
NC_001317.1	<i>Bacteriophage</i> 186	<i>Myoviridae</i>	P2virus	K	30624	53.09	46	Enterobacteriaceae	Escherichia
NC_005056.1	<i>Bacteriophage</i> WPhi	<i>Myoviridae</i>	P2virus	K	32684	51.72	44	Enterobacteriaceae	Escherichia
NC_022750.1	<i>Enterobacteria</i> phage fiAA91-ss	<i>Myoviridae</i>	P2virus	K	33628	51.91	40	Enterobacteriaceae	Escherichia
NC_028701.2	<i>Salmonella</i> phage SEN5	<i>Myoviridae</i>	P2virus	K	33509	53.36	47	Enterobacteriaceae	Salmonella
NC_029015.2	<i>Salmonella</i> phage SEN4	<i>Myoviridae</i>	P2virus	K	33509	53.36	47	Enterobacteriaceae	Salmonella
NC_021774.1	<i>Salmonella</i> phage FSL SP-004	<i>Myoviridae</i>	P2virus	K	29742	52.84	40	Enterobacteriaceae	Salmonella
NC_029003.2	<i>Salmonella</i> phage SEN1	<i>Myoviridae</i>	P2virus	K	29733	53.01	43	Enterobacteriaceae	Salmonella
NC_028943.1	<i>Escherichia</i> phage pro483	<i>Myoviridae</i>	P2virus	K	29237	52.98	43	Enterobacteriaceae	Escherichia
NC_019488.1	<i>Salmonella</i> phage RE-2010	<i>Myoviridae</i>	P2virus	K	34117	51.02	47	Enterobacteriaceae	Salmonella
NC_010463.1	<i>Enterobacteria</i> phage Fels-2	<i>Myoviridae</i>	P2virus	L	33693	52.49	46	Enterobacteriaceae	Salmonella
NC_026014.1	<i>Enterobacteria</i> phage P88	<i>Myoviridae</i>	P2virus	L	35814	52.87	53	Enterobacteriaceae	Escherichia
NC_019932.1	<i>Erwinia</i> phage ENT90	<i>Myoviridae</i>	P2virus	L	29564	55.81	60	Erwiniaceae	Erwinia
NC_010393.1	Phage Gifsy-2			M	45840	51.1	55	Enterobacteriaceae	Salmonella
NC_010392.1	Phage Gifsy-1			M	48491	51.1	58	Enterobacteriaceae	Salmonella
NC_001416.1	<i>Enterobacteria</i> phage lambda	<i>Siphoviridae</i>	Lambdavirus	M	48502	49.86	73	Enterobacteriaceae	Escherichia
NC_020845.1	Cyanophage MED4-213	<i>Myoviridae</i>	unclassified	N	180977	37.76	216	Prochloraceae	Prochlorococcus
NC_023693.1	<i>Enterobacteria</i> phage phi92	<i>Myoviridae</i>	unclassified	N	148612	37.43	250	Enterobacteriaceae	Escherichia
NC_009904.1	<i>Enterococcus</i> phage phiEF24C	<i>Myoviridae</i>	unclassified	N	142072	35.74	221	Enterococcaceae	Enterococcus
NC_001697.1	<i>Haemophilus</i> phage HP1	<i>Myoviridae</i>	Hp1virus	O	32355	40.01	42	Pasteurellaceae	Haemophilus
NC_003315.1	<i>Haemophilus</i> phage HP2	<i>Myoviridae</i>	Hp1virus	O	31508	39.94	37	Pasteurellaceae	Haemophilus
NC_005856.1	<i>Enterobacteria</i> phage P1	<i>Myoviridae</i>	P1virus	P	94800	47.31	110	Enterobacteriaceae	Escherichia
NC_031129.1	<i>Salmonella</i> phage SJ46	<i>Myoviridae</i>	P1virus	P	103445	48.58	122	Enterobacteriaceae	Salmonella
NC_010495.1	<i>Salmonella</i> phage E1	<i>Siphoviridae</i>	Pis4avirus	Q	45051	46.13	51	Enterobacteriaceae	Salmonella
NC_031924.1	<i>Salmonella</i> phage IME207	<i>Siphoviridae</i>	Pis4avirus	Q	47564	46.42	94	Enterobacteriaceae	Klebsiella
NC_005284.1	<i>Burkholderia</i> phage phi1026b	<i>Siphoviridae</i>	E125virus	NON	54865	60.68	83	Burkholderiaceae	Burkholderia
NC_024365.1	<i>Pseudomonas</i> phage phiPSA1	<i>Siphoviridae</i>	unclassified	NON	51090	58.57	51	Pseudomonadaceae	Pseudomonas
NC_031091.1	<i>Pseudomonas</i> phage MD8	<i>Siphoviridae</i>	unclassified	NON	43277	61.13	64	Pseudomonadaceae	Pseudomonas
NC_005859.1	<i>Enterobacteria</i> phage T5	<i>Siphoviridae</i>	T5virus	NON	121750	39.27	162	Enterobacteriaceae	Escherichia
NC_028748.2	<i>Bacillus</i> phage vB_BtS_BMBtp3	<i>Siphoviridae</i>	unclassified	NON	51366	35.45	76	Bacillaceae	Bacillus
NC_028841.1	<i>Bacteriophage</i> Lily	<i>Siphoviridae</i>	unclassified	NON	44952	42.73	74	Paenibacillaceae	Paenibacillus
NC_019401.1	<i>Cronobacter</i> phage vB_CsaM_GAP32	<i>Myoviridae</i>	unclassified	NON	358663	35.55	545	Enterobacteriaceae	Cronobacter
NC_009821.1	<i>Enterobacteria</i> phage Phi1	<i>Myoviridae</i>	Rb49virus	NON	164270	40.5	276	Enterobacteriaceae	Escherichia
NC_020079.1	<i>Escherichia</i> phage phAPEC8	<i>Myoviridae</i>	unclassified	NON	147737	39.15	269	Enterobacteriaceae	Escherichia
NC_004827.1	<i>Bacteriophage</i> Aphi23	<i>Myoviridae</i>	unclassified	NON	43033	42.46	66	Pasteurellaceae	Haemophilus
NC_019934.1	<i>Cronobacter</i> phage ENT39118	<i>Siphoviridae</i>	Hk97virus	NON	39012	53.06	38	Enterobacteriaceae	Cronobacter
NC_013594.1	<i>Escherichia</i> phage D108	<i>Myoviridae</i>	Muvirus	NON	37235	51.76	55	Enterobacteriaceae	Escherichia
NC_019455.1	<i>Haemophilus</i> phage SuMu	<i>Myoviridae</i>	Muvirus	NON	37151	41.87	55	Pasteurellaceae	Haemophilus
NC_028898.1	<i>Mannheimia</i> phage vB_MhM_587AP1	<i>Myoviridae</i>	P2virus	NON	35764	42.06	51	Pasteurellaceae	Mannheimia
NC_028896.1	<i>Escherichia</i> phage pro147	<i>Myoviridae</i>	P2virus	NON	32675	50.74	44	Enterobacteriaceae	Escherichia
NC_003313.1	<i>Vibrio</i> phage K139	<i>Myoviridae</i>	Hp1virus	NON	33106	48.9	44	Vibrionaceae	Vibrio
NC_019522.1	<i>Pectobacterium</i> phage ZF40	<i>Myoviridae</i>	unclassified	NON	48454	50.2	68	Pectobacteriaceae	Pectobacterium
NC_019927.1	<i>Cronobacter</i> phage ENT47670	<i>Myoviridae</i>	unclassified	NON	47611	51.59	46	Enterobacteriaceae	Cronobacter
NC_015295.1	<i>Erwinia</i> phage phiEt88	<i>Myoviridae</i>	unclassified	NON	47279	47.33	68	Erwiniaceae	Erwinia
NC_025458.1	<i>Shewanella</i> sp. phage 1_41	<i>Myoviridae</i>	unclassified	NON	43510	42.7	69	Shewanellaceae	Shewanella
NC_026611.1	<i>Edwardsiella</i> phage GF-2 DNA	<i>Myoviridae</i>	unclassified	NON	43129	51.27	82	Hafniaceae	Edwardsiella
NC_027995.1	<i>Escherichia</i> phage vB_EcoM_ECO1230-10	<i>Myoviridae</i>	Cvm10virus	NON	41666	53.37	56	Enterobacteriaceae	Escherichia
NC_028699.1	<i>Salmonella</i> phage SEN34	<i>Myoviridae</i>	unclassified	NON	40740	49.91	63	Enterobacteriaceae	Salmonella
NC_027339.1	<i>Enterobacteria</i> phage Sfi	<i>Myoviridae</i>	unclassified	NON	38389	50.12	65	Enterobacteriaceae	Shigella
NC_024369.2	<i>Vibrio</i> phage X29	<i>Myoviridae</i>	unclassified	NON	41569	46.07	67	Vibrionaceae	Vibrio
NC_019514.1	<i>Erwinia</i> phage vB_EamP-S6	<i>Podoviridae</i>	unclassified	NON	74669	52.09	115	Erwiniaceae	Erwinia
NC_025445.1	<i>Enterobacteria</i> phage J8-65	<i>Podoviridae</i>	unclassified	NON	40981	55.69	47	Enterobacteriaceae	Escherichia
NC_025443.1	<i>Salmonella</i> phage 9NA	<i>Podoviridae</i>	Nonanavirus	NON	52869	42.91	84	Enterobacteriaceae	Salmonella
NC_011551.1	<i>Bacteriophage</i> APSE-2	<i>Podoviridae</i>	unclassified	NON	39867	42.91	41	Enterobacteriaceae	Candidatus
NC_000935.1	<i>Acyrtosiphon pisum</i> bacteriophage APSE-1	<i>Podoviridae</i>	unassigned	NON	36524	43.89	54	Enterobacteriaceae	Candidatus
NC_009514.1	Phage cdtI DNA	<i>Siphoviridae</i>	unclassified	NON	47021	49.12	60	Enterobacteriaceae	Escherichia
NC_031264.1	<i>Brucella</i> phage BiPBO1	<i>Siphoviridae</i>	unclassified	NON	46877	53.32	86	Brucellaceae	Brucella
NC_001901.1	<i>Bacteriophage</i> N15	<i>Siphoviridae</i>	N15virus	NON	46375	51.17	60	Enterobacteriaceae	Escherichia
NC_005069.1	<i>Yersinia</i> phage PY54	<i>Siphoviridae</i>	unclassified	NON	46339	44.57	67	Yersiniaceae	Yersinia
NC_019716.1	<i>Enterobacteria</i> phage mEp460	<i>Siphoviridae</i>	unclassified	NON	44510	50.88	59	Enterobacteriaceae	Escherichia

NC_018843.1	<i>Salmonella</i> phage SSU5	<i>Siphoviridae</i>	unclassified	NON	103299	51.11	130	Enterobacteriaceae	Salmonella
NC_029028.1	<i>Enterobacteria</i> phage JenP1	<i>Siphoviridae</i>	Nonagvirus	NON	60754	43.23	87	Enterobacteriaceae	Escherichia
NC_028776.1	<i>Enterobacteria</i> phage CAjan	<i>Siphoviridae</i>	Seuratvirus	NON	59670	44.71	91	Enterobacteriaceae	Escherichia
NC_019545.1	<i>Salmonella</i> phage SPN3UB	<i>Siphoviridae</i>	unclassified	NON	47355	49.61	71	Enterobacteriaceae	Salmonella
NC_005857.1	<i>Klebsiella</i> phage phiKO2	<i>Siphoviridae</i>	unclassified	NON	51601	51.49	64	Enterobacteriaceae	Klebsiella
NC_016158.1	<i>Escherichia</i> phage HK639	<i>Siphoviridae</i>	unclassified	NON	49576	52.45	76	Enterobacteriaceae	Escherichia
NC_009552.2	<i>Geobacillus</i> virus E2	<i>Siphoviridae</i>	unclassified	NON	40863	44.79	71	Bacillaceae	Geobacillus
NC_018454.1	<i>Cronobacter</i> phage phiES15	<i>Siphoviridae</i>	unclassified	NON	39974	53.54	52	Enterobacteriaceae	Cronobacter
NC_015296.1	<i>Salmonella</i> phage Vi01	<i>Ackermannviridae</i>	Vi1virus	NON	157061	45.22	208	Enterobacteriaceae	Salmonella
NC_001609.1	<i>Enterobacteria</i> phage P4	<i>Unclassified</i> <i>Caudovirales</i>		NON	11624	49.53	14	Enterobacteriaceae	Escherichia
NC_023575.1	<i>Pseudomonas</i> phage vB_PaeP_Tr60_Ab31	<i>Unclassified dsDNA</i>		NON	45550	57.11	69	Pseudomonadaceae	Pseudomonas
NC_020850.1	<i>Vibrio</i> phage VBM1 genomic sequence	<i>Unclassified dsDNA</i>		NON	38374	42.26	56	Vibrionaceae	Vibrio
NC_010391.1	<i>Salmonella</i> phage Fels-1	<i>Unclassified bacterial viruses</i>		NON	42723	51.56	52	Enterobacteriaceae	Salmonella
NC_001954.1	<i>Enterobacteria</i> phage If1	<i>Inoviridae</i>	Escherichia virus If1	NON	8454	43.71	10	Enterobacteriaceae	Escherichia
NC_006294.1	<i>Vibrio</i> phage KSF-1phi	<i>Inoviridae</i>	Vibrio virus KSF1	NON	7107	44.38	12	Vibrionaceae	Vibrio
NC_001332.1	<i>Enterobacteria</i> phage I2-2	<i>Inoviridae</i>	Lineavirus	NON	6744	42.72	9	Enterobacteriaceae	Escherichia
NC_025824.1	<i>Enterobacteria</i> phage fd strain 478	<i>Inoviridae</i>	unclassified	NON	6408	40.89	10	Enterobacteriaceae	Escherichia

Table 2: Nucleotide identify matrix for 17 clusters of Cluster A-Q

Cluster A	F116	BPP1	BcepC6B
F116	100	45.1	51.5
BPP1		100	44.8
BcepC6B			100

Cluster B	KL3	RSY1	KS14	phiE255	BcepMu
KL3	100	42.5	63.1	45.6	44.3
RSY1		100	41.6	40.7	40.2
KS14			100	44.9	43.2
phiE255				100	84.3
BcepMu					100

Cluster C	D3112	DMS3
D3112	100	84.4
DMS3		100

Cluster D	epsilon34	118970_sal4	UAB_Phi20	ST104	SEN22	ST160	g341c	ST64T	103203_sal5	SPN9CC	vB_SemP_Emek	IME10	Sf6	Sf101	HK
epsilon34	100	61.7	71.0	60.8	70.6	57.9	85.9	58.4	70.00	62.1	64.1	59.5	53.2	61.1	54.
118970_sal4		100	73.1	82.1	52.5	78.3	58.4	76.6	98.8	92.7	58.8	59.8	49.0	48.6	54.
UAB_Phi20			100	79.5	65.9	76.6	66.6	76.4	89.0	71.7	66.1	51.4	48.7	53.7	50.
ST104				100	56.2	81.6	56.9	81.8	78.1	83.9	64.5	62.3	50.5	52.4	52.
SEN22					100	50.0	66.0	52.9	64.5	52.7	64.0	48.9	52.9	54.7	47.
ST160						100	54.0	86.9	69.2	74.8	58.2	63.4	49.1	51.1	57.
g341c							100	55.4	67.1	58.6	68.0	58.8	52.5	58.6	53.
ST64T								100	71.0	81.1	63.5	64.0	48.1	50.3	54.
103203_sal5									100	98.1	66.4	53.8	49.7	51.5	50.
SPN9CC										100	59.3	61.5	48.9	48.7	51.
vB_SemP_Emek											100	57.5	50.1	52.7	50.
IME10												100	70.1	62.7	57.
Sf6													100	55.01	50.
Sf101														100	53.
HK															100

Cluster E	TL-2011b	Tyrion	epsilon15	SPN1S
TL-2011b	100	65.5	58.8	59.2
Tyrion		100	61.8	70.0
epsilon15			100	73.2
SPN1S				100

Cluster F	VT2phi_272	Min27	Ss-VASD	POCJ13	933W	VT2-Sakai	P13374	Stx	86
VT2phi_272	100	80.5	90.2	74.5	81.7	81.5	92.0	78.0	88.5
Min27		100	79.0	75.9	96.8	94.0	81.7	79.3	85.8
Ss-VASD			100	76.8	80.1	80.3	94.0	80.2	85.7
POCJ13				100	76.6	77.0	75.8	93.7	84.1
933W					100	94.2	82.9	80.4	85.9
VT2-Sakai						100	83.1	80.0	85.6
P13374							100	80.0	88.0
Stx								100	83.8
86									100

Cluster G	BP-4795	YYZ-2008	1717	vB_SosS_Oslo	ES18
BP-4795	100	84.7	80.5	44.2	42.3
YYZ-2008		100	86.3	44.9	42.7
1717			100	44.7	42.5
vB_SosS_Oslo				100	65.1
ES18					100

Cluster H	mEp390	mEpX2	HK75	mEpX1	HK633	HK446	mEp235	HK022	HK97	HK106
mEp390	100	60.7	64.3	62.7	63.5	63.8	52.0	59.8	62.9	63.4
mEpX2		100	82.3	75.0	77.0	78.1	79.6	86.0	75.3	76.2
HK75			100	84.9	89.0	83.7	72.2	81.5	84	85.0
mEpX1				100	82.0	80.8	66.4	74.1	82.5	79.5
HK633					100	83.2	68.1	77.5	82.7	84.7
HK446						100	65.7	75.4	85.4	84.5
mEp235							100	83.8	66.8	68.7
HK022								100	77.6	78.0
HK97									100	88.0
HK106										100

Cluster I	phi80	HK225	mEp237	c-1
phi80	100	70.1	68.2	58.0
HK225		100	81.5	49.2
mEp237			100	51.2
c-1				100

Cluster J	118970_sal3	phiP27	SfII	ST64B	SfIV	SfV
118970_sal3	100	59.4	63.0	89.5	63.3	62.1
phiP27		100	62.1	68.2	65.5	57.0
SfII			100	61.2	81.4	83.1
ST64B				100	61.7	59.0
SfIV					100	79.5
SfV						100

Cluster K	P2	L413C	PsP3	186	WPhi	fiAA91-ss	SEN5	SEN4	SP-004	SEN1	pro483	RE-2010
P2	100	87.1	65.1	65.3	87.1	86.0	47.9	47.9	71.3	65.1	88.4	54.8
L413C		100	65.0	65.6	88.8	89.0	49.0	49.0	75.4	65.4	87.9	54.9
PsP3			100	79.2	65.4	64.4	47.5	47.6	79.5	91.8	70.8	56.8
186				100	65.2	64.8	48.1	48.1	76.3	80.2	71.3	57.0
WPhi					100	88.2	48.5	48.5	72.9	65.4	89.4	55.2
fiAA91-ss						100	48.1	48.1	73.4	64.7	89.0	55.4
SEN5							100	100	50.4	50.1	48.7	54.4
SEN4								100	50.4	50.1	48.7	54.4
SP-004									100	80.4	70.7	58.2
SEN1										100	71.2	60.0
pro483											100	51.3
RE-2010												100

Cluster L	Fels2	P88	ENT90
Fels2	100	41.3	52.7
P88		100	48.3
ENT90			100

Cluster M	Fels2	P88	ENT90
Fels2	100	41.3	52.7
P88		100	48.3
ENT90			100

Cluster N	MED4-213	phi92	phiEF24C
MED4-213	100	48.3	42.4
phi92	48.3	100	42.3
phiEF24C	42.4	42.3	100

Cluster O	HP1	HP2
HP1	100	90.4
HP2		100

Cluster P	P1	SJ46
P1	100	46.3
SJ46		100

Cluster Q	II-E1	IME207
II-E1	100	80.2
IME207		100

Table 3: The distribution of conserved Phamily members among clusters of *Salmonella* bacteriophages

Pham	Cluster (Number of presence)																
	A(3)	B(5)	C(2)	D(15)	E(4)	F(9)	G(5)	H(10)	I(4)	J(6)	K(12)	L(3)	M(3)	N(3)	O(2)	P(2)	Q(2)
6(10)						P											
27(17)								P	P								
35(14)													P				
45(9)									P								
53(4)													P				
58(3)			P														
64(4)	P																
89(5)							P										
103(2)	P																
124(3)																	P
127(7)	P				P												
163(11)						P											
172(9)						P											
180(12)				P													
195(7)														P			
212(3)			P														
239(12)						P											
269(23)											P	P			P		
297(10)										P							
312(10)						P											
316(9)						P											
321(7)					P												
329(3)			P														
333(4)			P														
375(5)	P																
393(13)				P													
403(3)																	P
413(2)																P	
415(5)					P												
447(17)				P													
450(5)							P										
450(2)																	P
460(6)																	P
474(15)																P	
475(9)										P							
489(4)	P																
519(19)								P		P							
520(9)						P											
526(12)				P													
528(8)										P							
529(10)						P											
550(9)						P											
573(19)											P						
617(19)		P									P						
640(26)				P													
669(9)										P							
708(9)										P							
728(2)														P			
735(19)											P	P					
738(20)											P	P					
746(7)										P							
769(19)											P				P		
776(4)															P		
779(11)																	P
788(18)											P						
795(13)				P													
804(3)															P		
824(7)	P				P												
825(9)						P											
890(6)							P										P
895(13)						P											
910(2)																	P
944(2)																P	
963(21)		P									P	P					
965(12)				P													
991(10)										P							
1026(15)													P				
1059(8)															P		
1115(2)																P	
1144(3)			P														
1171(9)						P											
1188(6)																	P
1190(18)								P									
1230(23)								P	P								

3082(2)									P
3125(21)	P					P	P		
3181(14)		P							
3246(12)							P		
3258(6)			P						
3330(2)									P
3331(2)									P
3381(3)								P	
3386(4)			P						
3508(23)				P					
3555(14)					P				
3597(5)			P						
3606(2)									P
3725(11)						P			
3826(17)				P					
3828(2)	P								
3975(4)			P						
4047(2)	P								
4053(13)				P					
4276(2)									P
4340(18)						P			
4411(3)	P								
4412(2)									P
4652(3)								P	
4878(4)			P						
4887(33)				P					
4959(25)		P							
5158(2)	P								
5196(2)								P	
5261(2)	P								
5722(2)	P								

Conserved families shown in bold are shared by at least two clusters and are not unique. Nevertheless, their presences contribute to the generation of cluster-specific profiles of families.

Figures

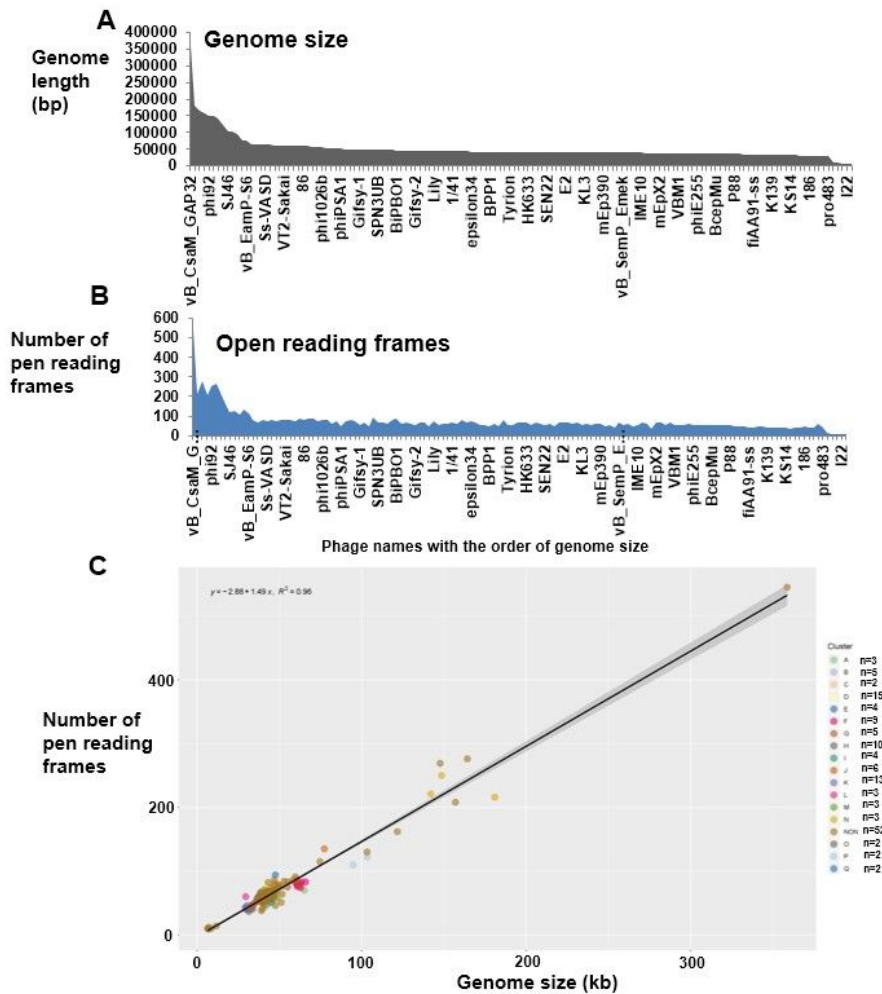
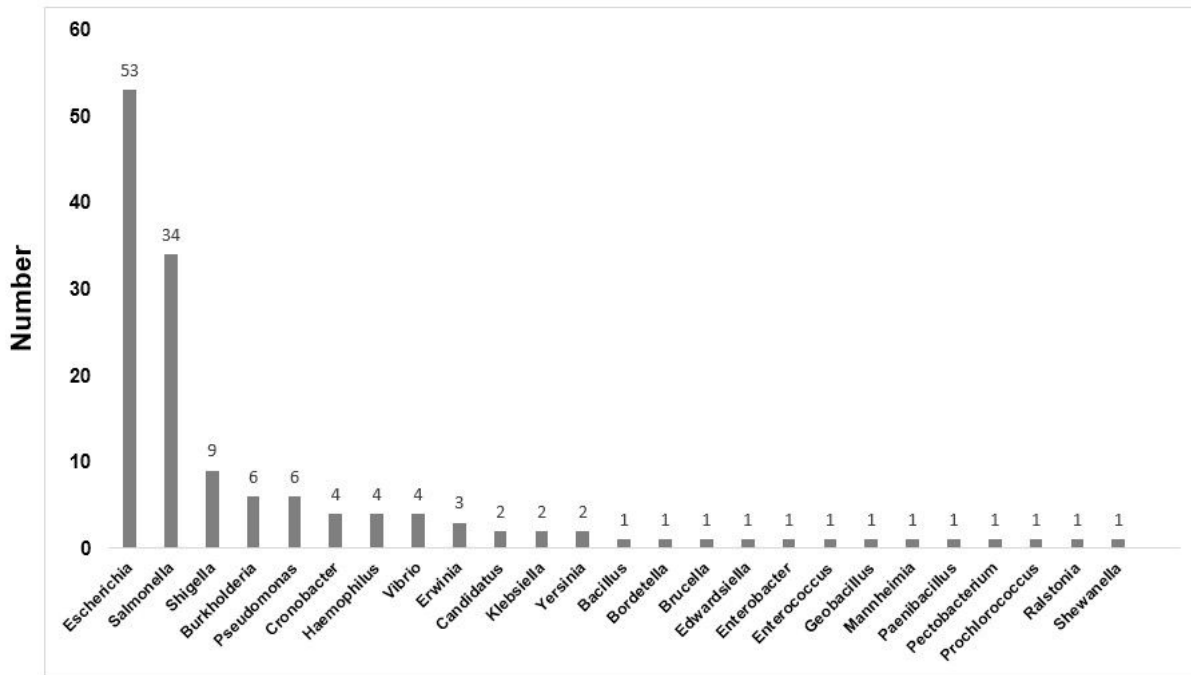


Figure 1

Genome characteristics of 142 Salmonella prophages. (A) Plot of genome sizes (B) Plot of the number of Open Reading Frames (ORFs). X axis shows names of each of the 142 prophages. Y axis represents either the genome length or number of detected ORFs in each prophage genome. (C) The correlation between the number of predicted ORFs and genome size in prophage genomes ($R^2 = 0.95$, $p < 0.001$). The shading besides the line indicates 95% confident interval of the linear correlation. The genomes from different clusters were shown with a different color of dot.



Bacterial hosts harbouring *Salmonella* prophages based on first description in the literature

Figure 2

Bacterial hosts of 142 *Salmonella* prophages. The X axis represents the number of prophages while the Y axis represents the frequency of occurrence in the bacterial host as identified in Virus-Host DB (<https://www.genome.jp/virushostdb/>).

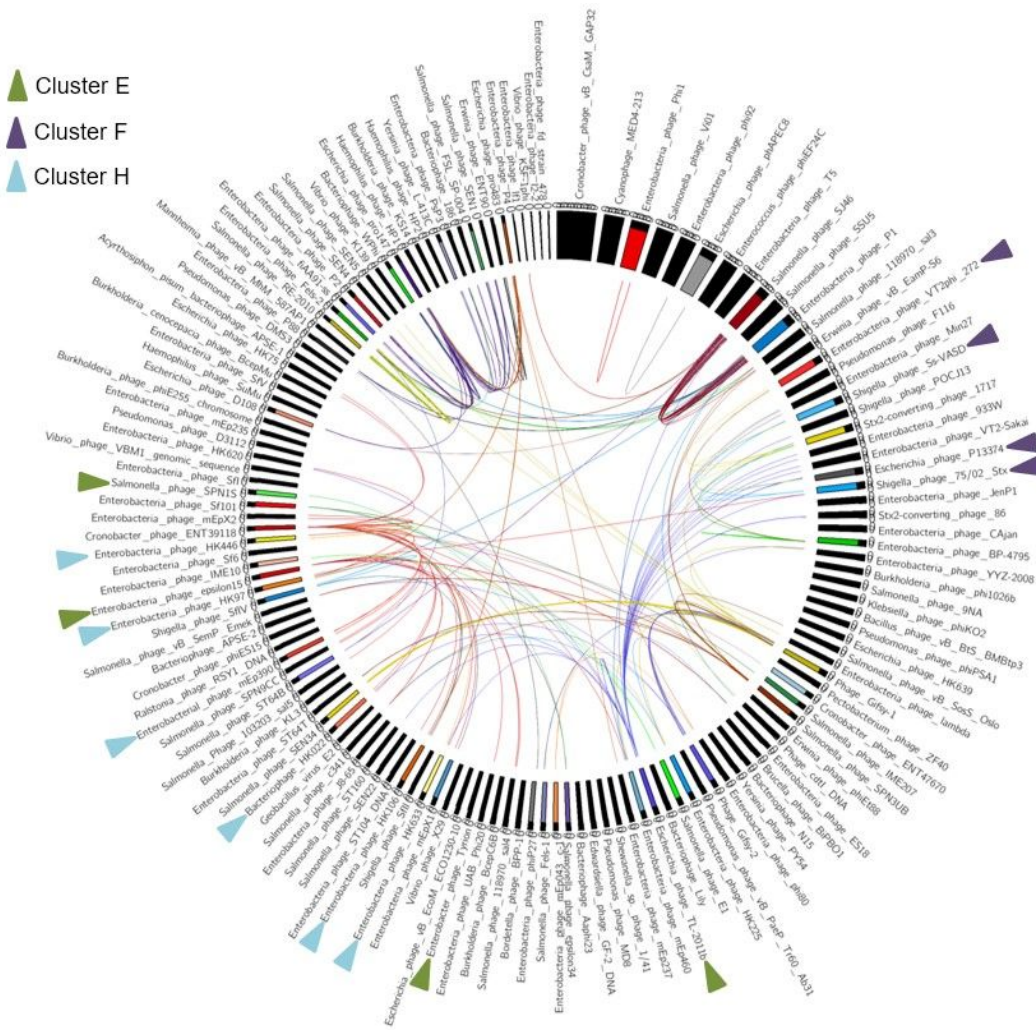


Figure 3

Similarities among 142 Salmonella prophages based on nucleotide identity and displayed using Circos. Nucleotide identities between prophages were calculated and coordinates were generated using NUCmer/MUMmer and displayed as Circos. Names of prophages are shown on the outer layer and arranged according to genome sizes. Prophages are highlighted in color block if more than one link (using the same color line as prophage block) existed with any of the other prophages. In contrast, prophages were shown in black block if no nucleotide similarity was detected with the other genomes.

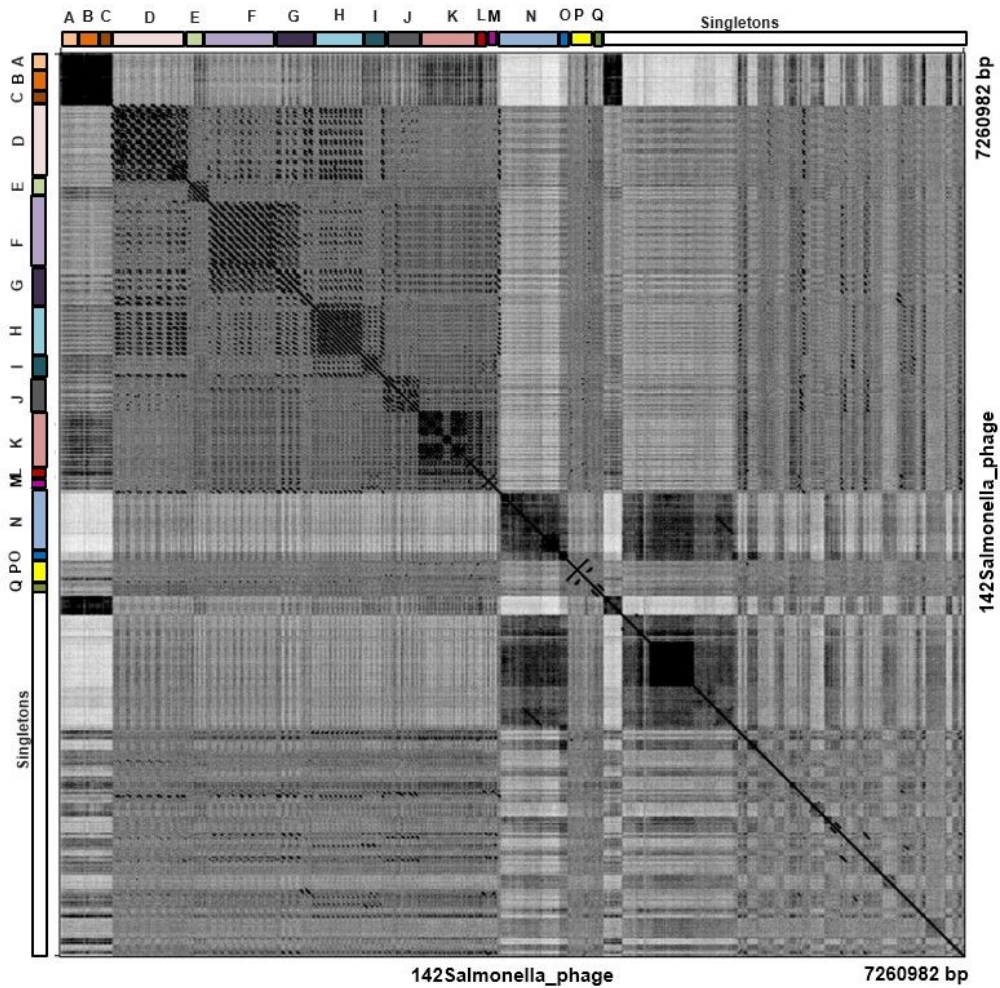


Figure 4
 Whole-genome dot plot comparison of prophage nucleotide sequences of Salmonella. Prophage genomes (n=142 phage) were concatenated into a single sequence with a total length of 7,260,982 bp, which plots against itself with a sliding window of 10 bp and visualized by Genome Pair Rapid Dotter (Gepard) 1.40 version. A total of 90 prophage genomes were assigned to 17 groups A - Q, and the remaining 52 prophage genomes plotted as singletons.

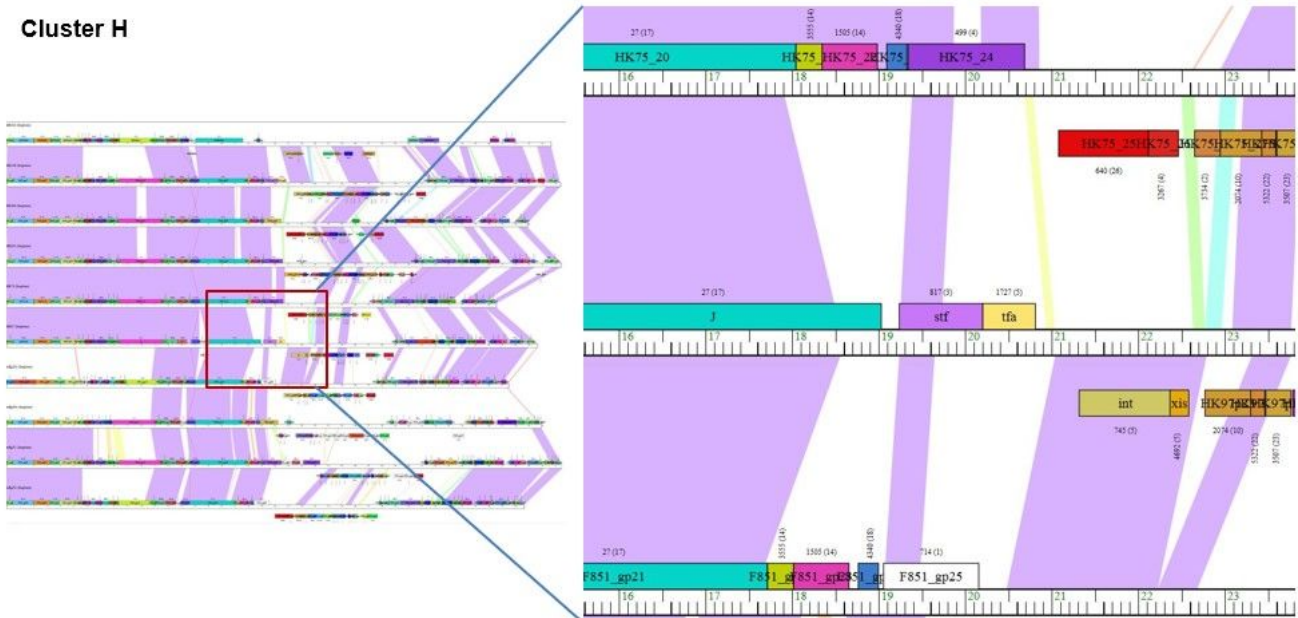


Figure 5

Genomic maps of Salmonella prophages belonging to Cluster H using the Phamerator software. (A) 11 prophage genomes (mEp390, mEpX2, HK75, mEpX1, HK633, HK446, mEp235, ENT39118, HK022, HK97 and HK106) was present in Cluster J. (B) Close-up view of partial of cluster J map. Blocks represent predicted ORFs, genes are color-coded according to their pham assignment. Gene names are shown within each gene box and the pham number and number of pham members are shown in parentheses above each gene. Shading between genomes indicates regions of pair-wise nucleotide similarity and was coded in color spectrum so that color indicates nucleotide similarity with violet being the most similar and red being the least similar. No shading suggests there is no similarity.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [TableS2DNAalignmentsfor142phagenucleotideidentity.xlsx](#)
- [TableS3AssignedPhamnumbersandcolorcode.xlsx](#)
- [FigureS1.Phylogenetictreeof142Salmonellaprophages.pptx](#)
- [TableS1Salmonellaentericaphageprofiles.xlsx](#)