

Insights into origin and evolution of α -proteobacterial gene transfer agents

Migun Shakya,^{1,†} Shannon M. Soucy,¹ and Olga Zhaxybayeva^{1,2,*,‡}

¹Department of Biological Sciences, Dartmouth College, 78 College Street, Hanover, NH 03755, USA and

²Department of Computer Science, Dartmouth College, 6211 Sudikoff Lab, Hanover, NH 03755, USA

*Corresponding author: E-mail: olgazh@dartmouth.edu

[†]Present address: Bioscience Division, P.O. Box 1663, Los Alamos National Laboratory, Los Alamos, NM 87544, USA.

[‡]<http://orcid.org/0000-0002-1809-3909>

Abstract

Several bacterial and archaeal lineages produce nanostructures that morphologically resemble small tailed viruses, but, unlike most viruses, contain apparently random pieces of the host genome. Since these elements can deliver the packaged DNA to other cells, they were dubbed gene transfer agents (GTAs). Because many genes involved in GTA production have viral homologs, it has been hypothesized that the GTA ancestor was a virus. Whether GTAs represent an atypical virus, a defective virus, or a virus co-opted by the prokaryotes for some function, remains to be elucidated. To evaluate these possibilities, we examined the distribution and evolutionary histories of genes that encode a GTA in the α -proteobacterium *Rhodobacter capsulatus* (RcGTA). We report that although homologs of many individual RcGTA genes are abundant across bacteria and their viruses, RcGTA-like genomes are mainly found in one subclade of α -proteobacteria. When compared with the viral homologs, genes of the RcGTA-like genomes evolve significantly slower, and do not have higher %A+T nucleotides than their host chromosomes. Moreover, they appear to reside in stable regions of the bacterial chromosomes that are generally conserved across taxonomic orders. These findings argue against RcGTA being an atypical or a defective virus. Our phylogenetic analyses suggest that RcGTA ancestor likely originated in the lineage that gave rise to contemporary α -proteobacterial orders Rhizobiales, Rhodobacterales, Caulobacterales, Parvularculales, and Sphingomonadales, and since that time the RcGTA-like element has co-evolved with its host chromosomes. Such evolutionary history is compatible with maintenance of these elements by bacteria due to some selective advantage. As for many other prokaryotic traits, horizontal gene transfer played a substantial role in the evolution of RcGTA-like elements, not only in shaping its genome components within the orders, but also in occasional dissemination of RcGTA-like regions across the orders and even to different bacterial phyla.

Key words: exaptation; domestication; horizontal gene transfer; bacterium-virus co-evolution; bacteriophage.

1. Introduction

Prokaryotes are hosts not only of their own genetic material, but also of mobile genetic elements, a broad class of entities that includes integrated viruses (prophages) (Frost et al. 2005). Traditionally, viruses are viewed as selfish genetic elements, but in some instances they can serve as vectors of horizontal

gene transfer (HGT) (Touchon et al. 2017), an important driver of evolutionary success of many cellular lifeforms (Zhaxybayeva and Doolittle 2011; Koonin 2016). In other instances, prophage-like elements can provide immunity against other viruses (Canchaya et al. 2003), and therefore are beneficial to their host. While the evolutionary relationship between viruses and cellular lifeforms is still intensely debated, co-option of genes by

both empires is likely frequent (Krupovic and Koonin 2017). For example, a few prokaryotic cellular functions that are associated with adaptations, such as cell–cell warfare, microbe–animal interactions, HGT, and response to environmental stress, are carried out using structures that resemble viral ‘heads’ (e.g. Sutter et al. 2008; McHugh et al. 2014) or ‘tails’ (e.g. Shikuma et al. 2014; Borgeaud et al. 2015). These structures appear to be widespread across archaea and bacteria (e.g. Sarris et al. 2014; Böck et al. 2017). In other instances, viral components appear to have originated from the cellular proteins (Krupovic and Koonin 2017). Such intertwined history of genes from cellular and non-cellular lifeforms resulted in a spectrum of genetic elements that, from the cellular host ‘point of view’, span from parasitic to benign to beneficial.

One class of such genetic elements is yet to be placed in this spectrum. Several unrelated bacterial and archaeal lineages are observed to produce particles that morphologically resemble viruses, yet appear to carry random pieces of host DNA instead of their own viral genome (Marrs 1974; Rapp and Wall 1987; Humphrey et al. 1997; Bertani 1999; Berglund et al. 2009). These particles were shown to transfer DNA between cells through a transduction-like mechanism (McDaniel et al. 2010; Lang et al. 2012), and that’s why they were dubbed ‘gene transfer agents’ (GTAs) (Marrs 1974). Since the transferred DNA may benefit the recipient cells, GTAs were postulated to be viruses that were ‘domesticated’ by prokaryotes (Bobay et al. 2014) and maintained by them as a mechanism of HGT (Lang et al. 2012). Yet, the details of such co-option event as well as impact of the GTA-mediated HGT remain to be deciphered, and alternatives, such as GTA being instead a defective or atypical virus, need to be evaluated.

The laboratory studies of *Rhodobacter capsulatus*’ GTA (or RcGTA for short) show that this element neither has a typical genome of a bacterial virus nor acts like a typical lysogenic virus. Its genome is distributed among five loci dispersed across the *R. capsulatus* chromosome (Hynes et al. 2016 and Fig. 1A). Seventeen of the genes are found in a single locus that has a genomic architecture typical of a siphovirus (Lang and Beatty 2007). Most of the genes in the locus encode proteins involved in head and tail morphogenesis of the RcGTA particle and are thus referred to as the ‘head–tail’ cluster (after Lang et al. 2017). The remaining four loci contain seven genes shown to be important for RcGTA production, release, and DNA transfer (Fogg et al. 2012; Hynes et al. 2012, 2016; Westbye et al. 2015). Even if RcGTA would preferentially package the regions of DNA that correspond to its own genome, the small head size of RcGTA particle can accommodate only approximately one-fifth of the genome (Lang et al. 2017), and therefore RcGTA can propagate itself only with the division of the host cell. Expression of the RcGTA genes, as well as production and release of RcGTA particles can be triggered by phosphate concentration (Leung et al. 2010; Westbye et al. 2013), salinity (McDaniel et al. 2012), and quorum sensing (Schaefer et al. 2002; Brimacombe et al. 2013). The latter is regulated via CckA–ChpT–CtrA phosphorelay (Leung et al. 2012; Mercer and Lang 2014), a widely used bacterial signaling network that also controls cell cycle (Chen and Stephens 2007; Mann et al. 2016) and flagellar motility (Zan et al. 2013). These observations suggest that RcGTA is under control of the bacterial host and is well-integrated with the host’s cellular systems. When did such integration occur? Could RcGTA still represent a selfish genetic element? Through phylogenomic analyses of a much more extensive genomic data set, we show that this α -proteobacterial GTA is a virus-related element whose genome evolves differently from what is expected

of a typical bacterial virus. We also infer that RcGTA-like element likely originated in an ancestor of an α -proteobacterial subclade that gave rise to at least five contemporary taxonomic orders. Since that time, the genomes of these elements were shaped by both co-evolution with the bacterial hosts and HGT.

2. Methods

2.1 Detection of RcGTA homologs in viruses and bacteria

Twenty-four genes from the RcGTA genome and their homologs from *Rhodobacterales* were used as queries in BLASTP (E-value <0.001; query and subject overlap by at least 60% of their length) and PSI-BLAST searches (E-value <0.001; query and subject overlap by at least 40% of their length; maximum of six iterations) against viral and bacterial databases using BLAST v. 2.2.30+ (Altschul et al. 1997). The viral database consisted of 1,783 genomes of dsDNA bacterial viruses extracted from the RefSeq database (release 76; last accessed on 1 July 2016) using ‘taxid = 35237 AND host = bacteria’ as a query. Bacterial RefSeq database (release 76; last accessed on 29 June 2016) and 255 completely sequenced α -proteobacterial genomes (Supplementary Table S1; genomes downloaded from ftp://ftp.ncbi.nlm.nih.gov/genomes/Bacteria/ on 1 November 2015) served as two bacterial databases. The detected RefSeq homologs were mapped to their source genomes using ‘release76.AutonomousProtein2Genomic.gz’ file (ftp.ncbi.nlm.nih.gov/refseq/release/release-catalog/, last accessed on 29 June 2016).

2.2 Assignment of genes to viral gene families

All annotated protein-coding genes from 255 α -proteobacterial genomes were used as queries in BLASTP searches (E-value <0.001; query and subject overlap by at least 40% of their length) against the Phage Orthologous Group (POG) database (Kristensen et al. 2013), which was downloaded from ftp://ftp.ncbi.nlm.nih.gov/pub/kristensen/thousandgenomespogs/blastdb/ in February 2016. The genes were assigned to a POG family based on the POG affiliation of the top-scoring BLASTP match.

2.3 Identification of RcGTA homologs adjacency in the bacterial genomes

Relative positions of RcGTA homologs in each genome were calculated from genome annotations in the GenBank feature tables downloaded from ftp.ncbi.nlm.nih.gov/genomes/all/ on 29 June 2016. Initially, a region of the genome was classified as putative RcGTA-like region if it had at least one RcGTA homolog. The identified regions were merged if two adjacent regions were separated by <15 open reading frames (ORFs). If the resulting merged region contained at least nine RcGTA homologs, it was designated as a large cluster (LC). The remaining regions were labeled as small clusters (SCs).

2.4 Examination of viral characteristics of RcGTA-like genes and regions

Prophage-like regions in 255 α -proteobacterial genomes were detected using PhiSpy v. 2.3 (Akhter et al. 2012) with ‘-n=5’ (at least five consecutive prophage genes), ‘-w=30’ (a window size of 30 ORFs), and ‘-t = 0’ (Generic Test Set) settings. An RcGTA-like region was classified as a putative prophage if it overlapped with a PhiSpy-detected prophage region.

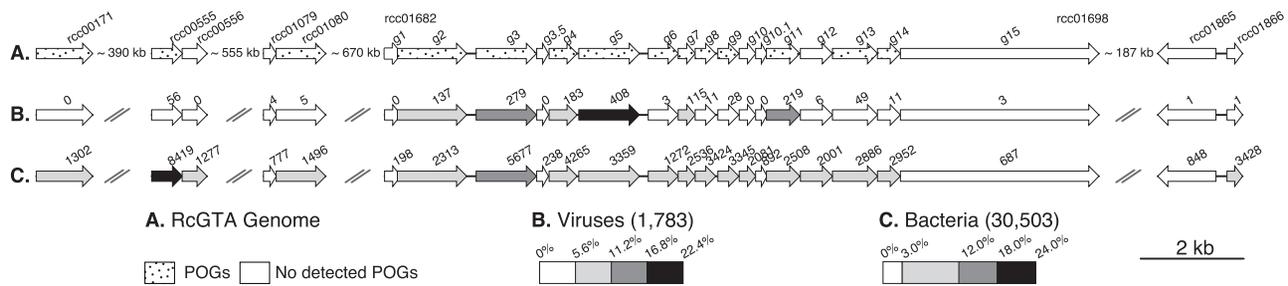


Figure 1. RcGTA genome and distribution of its homologs in viral and bacterial genomes. (A) The RcGTA genome architecture. The RcGTA genome consists of twenty-four genes scattered across five loci in the *R. capsulatus* SB1003 genome. The genes are represented by arrows. The majority of the ‘head–tail’ cluster genes (*R. capsulatus* SB1003 locus tags *rcc01682*–*rcc01698*; also known as *g1*–*g15*) encode genes involved in head and tail morphogenesis, *rcc00171* (*tsp*)—a tail spike protein (Hynes et al. 2016), *rcc00555* and *rcc00556*—endolysin and holin (Fogg et al. 2012; Hynes et al. 2012; Westbye et al. 2013), *rcc01079* (*ghsA*) and *rcc01080* (*ghsB*)—a head spike (Westbye et al. 2015), and *rcc01865* and *rcc01866*—putative regulatory elements (Hynes et al. 2016). The GenBank functional annotations of the genes are listed in Supplementary Table S4. Genes with similarity to gene families in the POG database (Kristensen et al. 2013) are shaded and their assigned POG categories are listed in Supplementary Table S4. (B) Presence of RcGTA gene homologs in 1,783 viral genomes. (C) Presence of RcGTA gene homologs in 30,503 bacterial genomes. The numbers above arrows represent the total number of detected homologs for each RcGTA gene, while the percent of viral (B) bacterial (C) genomes in which these homologs are found is depicted in shades of gray (see figure inset for scale). Genomic regions were visualized using R package *genoplR* (Guy et al. 2010).

Presence of viral integrases was assessed within all prophage-like and RcGTA-like regions, as well as among five ORFs upstream and downstream of the regions. The integrase homologs were identified using the *hmmsearch* program of the HMMER package v. 3.1b2 (Eddy 1998; E-value $<10^{-5}$), with HMM profiles of integrase genes from the pVOG database 2016 update (Grazziotin et al. 2017) as queries (VOG# 221, 275, 375, 944, 2142, 2405, 2773, 3344, 3995, 4650, 5717, 6225, 6237, 6282, 6466, 7017, 7518, 8218, 8244, and 10948) and 255 α -proteobacterial genomes as the database.

For both LCs and SCs, %G + C of the corresponding nucleotide sequence (defined as DNA sequence between the first and last nucleotide of the first and last RcGTA homolog of the region, respectively) was calculated using in-house Python scripts. The relative difference in GC content of the RcGTA-like region and the whole host genome was calculated by taking the difference in their %G + C and dividing it by the %G + C of the host genome. This correction allowed us to compare changes in the GC content of the RcGTA-like regions among genomes with variable GC contents. To account for possible heterogeneity of GC content within each genome, 100 genomic regions of the same size were randomly selected from the genome, their relative GC content calculated as described above and compared with that of the RcGTA-like region(s) using pairwise *t*-test.

Pairwise phylogenetic distances (PPD) among RcGTA homologs in viruses, LCs and SCs were calculated in RAXML v. 8.1.3 (Stamatakis 2014) using ‘-f x’ and ‘-m PROTGAMMAAUTO’ options. The latter choice selects the best amino acid substitution model and uses Γ distribution with four discrete rate categories to correct for among site rate variation (Yang 1994).

2.5 Identification of gene families

Sequences of protein-coding genes in 255 α -proteobacterial genomes were subjected to all-against-all BLASTP searches (v. 2.2.30+, E-value <0.0001), in which reciprocal top-scoring BLASTP matches were retained. The bit scores of the matches were converted to pairwise distances, which were used to form gene families via Markov clustering with the inflation parameter set to 1.2 (van Dongen 2000), as implemented in OrthoMCL v. 2.0.9 (Li et al. 2003). The gene presence/absence in the resulting 33,048 gene families was used to as a proxy of gene family conservation within and across α -proteobacteria.

2.6 Characterization of immediate gene neighborhoods of SCs and LCs

ORFs without detectable similarity to RcGTA genes, but found immediately upstream, downstream, or within ‘head–tail’ cluster, were classified based on (1) conservation within α -proteobacterial genomes and (2) potential viral origin. The proportion of α -proteobacterial genomes that have an ORF homolog was used as a proxy for conservation. The ORF was designated as viral if it was assigned to a POG.

2.7 Reference phylogeny of α -proteobacteria

From the list of 104 gene families used to reconstruct α -proteobacterial phylogeny (Williams et al. 2007), 99 gene families that are present in 80% of the 255 α -proteobacterial genomes were selected (Supplementary Table S2). For each gene family, homologs from *Geobacter sulfurreducens* PCA, *Ralstonia solanacearum* GMI1000, *Chromobacterium violaceum* ATCC 12472, *Xanthomonas axonopodis* pv. *citri* str. 306, *Pseudomonas aeruginosa* PAO1, *Vibrio vulnificus* YJ016, *Pasteurella multocida* subsp. *multocida* str. Pm70, *Escherichia coli* str. K12 substr. DH10B were added as an outgroup. The amino acid sequences of the resulting gene set were aligned using MUSCLE v. 3.8.31 (Edgar 2004). The best substitution model (listed in Supplementary Table S2) was selected using *ProteinModelSelection.pl* script downloaded from <http://sco.h-its.org/exelixis/resource/download/software/ProteinModelSelection.pl> in August 2016. Among site rate variation was modeled using Γ distribution with four rate categories (Yang 1994). All gene sets were combined into a single dataset with ninety-nine data partitions (one per gene). The maximum likelihood phylogenetic tree was reconstructed in RAXML v. 8.1.3 (Stamatakis 2014), using the best substitution model for each partition and twenty independent tree space searches. One hundred bootstrap samples were analyzed as described above, and the bootstrap support values $>60\%$ were mapped to the maximum likelihood tree. Reference phylogeny of only LC-containing taxa was obtained by pruning taxa that only have SCs.

2.8 Phylogenetic analysis of the large-cluster gene families

Since phylogenetic histories of many RcGTA homologs are poorly resolved in *Rhodobacterales* (Hynes et al. 2016), amino acid sequences of LC gene families were combined into one

'LC-locus' dataset in order to increase the number of phylogenetically informative sites. To ensure that such concatenation does not result in a mixture of genes with different evolutionary histories, phylogenetic trees reconstructed from the alignments of individual LC genes and of the concatenated LC-locus were compared. Amino acid sequences of gene families consisting only of LC homologs from α -proteobacterial genomes were aligned using MUSCLE v. 3.8.31 (Edgar 2004). The individual gene family alignments were concatenated into the 'LC-locus' alignment. Since many LCs do not have homologs of all RcgTA 'head-tail' genes, in the concatenated alignment the absences were designated as missing data (Stamatakis 2014). Additionally, 'taxa-matched' LC-locus alignments (i.e. concatenated alignments pruned to contain taxa found only in a specific LC gene family) were created. Maximum likelihood trees were reconstructed from the concatenated and individual LC gene alignments, as well as from the taxa-matched LC-locus alignments, using RAxML v. 8.2.9 (Stamatakis 2014) under the LG + Γ substitution model (Yang 1994; Le and Gascuel 2008). The LG model was determined as the best substitution model for each individual gene set using the PROTGAMMAAUTO option of RAxML v. 8.2.9 (Stamatakis 2014). For each LC gene, consensus trees were reconstructed from ten bootstrap samples of taxa-matched LC-locus alignments. These consensus trees were compared with 100 bootstrap sample trees for the corresponding LC gene. The congruence between the phylogenies was measured using relative 'Tree Certainty All' (TCA) values (Salichos et al. 2014), as implemented in RAxML v. 8.2.9 (Stamatakis 2014). The TCA values were classified into 'strongly conflicting' ($-1.0 \leq \text{TCA} \leq -0.7$), 'moderately conflicting' ($-0.7 < \text{TCA} < 0.7$), and 'not conflicting' ($0.7 \leq \text{TCA} \leq 1.0$) categories. None of the individual LC gene phylogenies had strongly conflicting TCA values, although six out of the seventeen LC gene phylogenies had moderately conflicting values (Supplementary Table S3). Since an uncertain position of just one taxon can affect support values of multiple bipartitions (Aberer et al. 2013) and TCA values calculated from them, a site specific placement bias (SSPB) analysis (Berger et al. 2011) was carried out, as implemented in RAxML v. 8.2.9 (Stamatakis 2014). The LC-locus phylogeny and the LC-locus alignment were used as the reference tree and the input sequence alignment, respectively. The sliding window size was set to one hundred amino acids. The SSPB analysis revealed that for fifteen of the seventeen genes only three nodes, on average, separate the optimal positions of each taxon in the gene and LC-locus phylogenies, indicating largely compatible evolutionary histories of individual LC genes and of their combination (Supplementary Fig. S1). Consequently, we decided to use the combined, taxa-matched LC-locus alignment for further phylogenetic analyses.

To assess the similarity of the evolutionary histories of LC loci and of their hosts, the topologies of the reference α -proteobacterial tree were compared with one hundred bootstrap sample trees reconstructed from the taxa-matched LC-locus alignment assembled as described above. The congruence was quantified using Internode Certainty (IC) values (Salichos et al. 2014), as implemented in RAxML v.8.2.9 (Stamatakis 2014). The IC values were classified into 'strongly conflicting' ($-1.0 \leq \text{IC} \leq -0.7$), 'moderately conflicting' ($-0.7 < \text{IC} < 0.7$), and 'not conflicting' ($0.7 \leq \text{IC} \leq 1.0$) categories.

2.9 Conservation of gene neighborhoods across phylogenetic distance

Forty ORFs upstream and downstream of the ATP synthase operon, the ribosomal protein operon, and the LCs were

extracted from eighty-seven α -proteobacterial genomes with at least one LC. The same procedure was performed for a transposase from IS3/IS911 family, which was detected in eighteen of the eighty-seven genomes. Each ORF was assigned to a gene family, and the pairwise conservation of the corresponding flanking regions was calculated as a proportion of gene families shared between a pair of genomes. PPDs of 16S rRNA genes were used as a proxy for time t . The 16S rRNA genes were aligned against the Greengenes database (Desantis et al. 2006; last accessed in August 2016) using mothur v. 1.35.1 (Schloss et al. 2009). The PPDs were calculated from the alignment using RAxML v. 8.1.3 (Stamatakis 2014) under GTR + Γ substitution model. Based on Model 3 from Rocha (2006), the gene neighborhood decay was modeled as p^t , where p is the probability of the 40 ORFs to remain in the same genomic region and t is time. This model assumes that every gene has equal probability of separating from the region. The parameter p was estimated, and the fit between the model and data was assessed using nonlinear least squares method (Bates and Watts 1988), as implemented in the *nls* R function.

2.10 Examination of genes that flank RcgTA-like clusters in non- α -proteobacterial genomes

Amino acid sequences of five ORFs upstream and downstream of the RcgTA-like clusters in the genomes of actinobacteria *Streptomyces purpurogeniscleroticus* NRRL B-2952 (GenBank accession number LGEI00000000.1; Ju et al. 2015) and *Asanoa ferruginea* NRRL B-16430 (LGEJ00000000.1; Ju et al. 2015), a γ -proteobacterium *Pseudomonas bauzanensis* W13Z2 (JFHS00000000.1; Wang et al. 2014), and a cyanobacterium *Scytonema millei* VB511283 (JTJC00000000.1; Sen et al. 2015) were searched against the database of protein-coding genes of 255 α -proteobacterial genomes using BLASTP (E-value < 0.001 ; query and subject overlap by at least 60% of their length). Three genes that flank LCs in *P. bauzanensis*, *S. millei*, and a few α -proteobacteria (*kpdE*, *pbpC*, *dnaJ*) were used as queries to detect their homologs in other cyanobacteria and γ -proteobacteria in BLASTP searches of γ -proteobacteria and cyanobacteria in the *nr* database (<https://blast.ncbi.nlm.nih.gov/Blast.cgi>, last accessed on 31 October 2016; E-value < 0.001 ; query and subject overlap by at least 60% of their length). For each of the three gene sets, four homologs from each identified taxonomic group were selected and combined with all available homologs from LC-containing α -proteobacteria. The amino acid sequences of the combined datasets were aligned using MUSCLE v. 3.8.31 (Edgar 2004). Maximum likelihood trees were reconstructed using FastTree v. 2.1.8 (Price et al. 2010) under JTT + Γ substitution model. The trees were visualized and annotated using EvolView (He et al. 2016; <http://evolgenius.info/evolview>, last accessed on 13 November 2016).

2.11 Identification of CRISPR/Cas defense systems

Spacers of the clustered regularly interspaced short palindromic repeats (CRISPRs) were identified by scanning 255 α -proteobacterial genomes for both CRISPR leader sequences and repeat structures using CRISPRleader v. 1.0.2 (Alkhnbashi et al. 2016) and PILER-CR v. 1.06 (repeat length between sixteen and sixty-four; spacer length between eight and sixty-four; minimum array length of three repeats; minimum conservation of repeats 0.9; minimum repeat length ratio 0.9; minimum spacer length ratio 0.75) (Edgar 2007), respectively. CRISPR-associated proteins were identified using the *hmmsearch* program of the HMMER package v. 3.1b2 (Eddy 1998; E-value $< 10^{-5}$), with TIGRFAM

families of the CRISPR-associated proteins as queries (TIGRFAM # 287, 372, 1573, 1587, 1596, 1863, 1865, 1868, 1869, 1876, 1907, 2562, 2589, 2590, 2593, 2621, 3158, 3637, 3638, 3639, 3640, 3641, 3983; Selengut et al. 2007), and protein-coding genes of 255 α -proteobacterial genomes as database. Genomes were designated to have putative functional CRISPR/Cas systems if both CRISPR spacer arrays and at least three CRISPR-associated proteins from the same CRISPR class (Makarova et al. 2015) were present. Spacer sequences from these putatively functional CRISPR systems were used as BLASTN (v. 2.4.0+; E-value <10; task = tblastn-short) queries against all gene sequences in LCs.

2.12 Identification of decaying RcGTA ‘head–tail’ cluster homologs

The nucleotide sequences of LC and SC regions were extracted from each genome and translated in their entirety into all six reading frames. RcGTA ‘head–tail’ cluster genes were used as BLASTP (v. 2.4.0+, E-value ≤ 0.001) queries in searches against the database of translated LC and SC regions. Homologous sequences within LC and SC regions that did not match the annotated ORFs were designated as putatively decaying RcGTA genes.

3. Results

3.1 Many RcGTA genes share evolutionary history with viruses

Consistent with the earlier proposed hypothesis that RcGTA is an element originated from a virus (Lang and Beatty 2000), fourteen out of the twenty-four RcGTA genes can be assigned to a POG (Fig. 1A and Supplementary Table S4), and eighteen out of the twenty-four genes have at least one readily detectable homolog in available genomes of *bona fide* bacterial viruses (Fig. 1B). However, none of the bacterial viruses in RefSeq database contained a substantial number of the RcGTA homologs (at most eight of them in *Rhizobium* phage 16-3 [GenBank Acc. No. DQ500118.1]), and no α -proteobacterial genomes with CRISPR systems contained significant matches between CRISPR spacers and RcGTA genome, suggesting that the presumed progenitor virus is either extinct or remains unsampled. Collectively, RcGTA-like genes are unevenly represented across viral genomes (Fig. 1B). Among the five most abundant are homologs of *rcc001683* (*g2*), *rcc001684* (*g3*), *rcc001686* (*g4*), and *rcc001687* (*g5*), which encode terminase, portal protein, prohead protease, and major capsid protein, respectively. Their widespread occurrence is not surprising. First, these genes correspond to the so-called ‘viral hallmark genes’, defined as genome replication and virion formation genes with a wide distribution among diverse viral genomes (Koonin et al. 2006; Iranzo et al. 2016). Second, *g3*, *g4*, and *g5* belong to HK97 family, and HK97-like genes are common in Caudovirales (Iranzo et al. 2017)—one of the most abundant viral groups (Cobián Güemes et al. 2016). However, thirteen of the RcGTA genes are either sparsely represented (found in ten or fewer viral genomes; seven genes) or not detected at all (six genes) in viruses. These genes may be fast-evolving viral genes, auxiliary (non-hallmark) viral genes, which could be specific to the viral lineage that gave rise to RcGTA, non-viral (cellular) genes, or simply be a result of undetected homology due to the paucity of viral genomes in GenBank. The latter possibility cannot be ignored, especially given that (1) three of the thirteen genes are assigned to a POG (Fig. 1A), and therefore are likely of viral origin, and (2) only 94 of the 1,783 screened bacterial viruses (5%) have α -proteobacteria as

their known host. Since about half of the bacterial genomes are estimated to contain at least one prophage in their genome (Touchon et al. 2016), inclusion of RcGTA homologs from bacterial genomes could reduce the possibility of undetected homology due to lack of available viral genomes. Therefore, we looked for RcGTA homologs in all bacterial RefSeq records, and in 255 completely sequenced α -proteobacterial genomes in particular (Supplementary Table S1).

3.2 Homologs of RcGTA genes are abundant in bacteria, but RcGTA-like genomes are mainly found in class α -proteobacteria

Every RcGTA gene has at least 198 homologs in RefSeq database (Fig. 1C), and therefore even the 13 ‘rare’ genes are relatively abundant in bacterial genomes. Moreover, 6 of the 13 genes (*rcc00171*, *rcc00556*, *rcc01688* [*g6*], *rcc01692* [*g10*], *rcc001695* [*g12*], and *rcc01866*), including 3 assigned to POGs, are found in more than 1,000 bacterial genomes (Fig. 1C). Therefore, it is unlikely that the thirteen rare genes are auxiliary viral genes specific to the viral lineage that gave rise to RcGTA. However, a larger number of genomes of α -proteobacterial viruses are needed to identify if these genes are of viral or cellular origin.

Collectively, homologs of RcGTA genes are found within thirteen bacterial phyla (Fig. 1C and Supplementary Fig. S2a). Interestingly, 84 and 35% of the identified homologs belong to the phylum Proteobacteria and to its class α -proteobacteria, respectively. Furthermore, many α -proteobacterial genomes contain more ‘head–tail’ RcGTA homologs than any available viral genome, and these homologs are often clustered (Supplementary Fig. S2b). In contrast, most other taxa have only few RcGTA homologs that are scattered across the genome (Supplementary Fig. S2b and c). Since in the RcGTA genome the ‘head–tail’ genes are in one locus (Fig. 1A) and are presumed to be co-transcribed (Kuchinski et al. 2016), functional GTAs of this family might be restricted mainly to α -proteobacteria.

Within α -proteobacteria, homologs of RcGTA genes are detectable in 197 out of the 255 examined genomes (77%). These homologs are distributed across eight α -proteobacterial orders with at least one available genome and two unclassified α -proteobacterial genera (Fig. 2). While thirteen out of the seventeen ‘head–tail’ genes are abundant across α -proteobacteria (and across Bacteria in general), only two of the genes from the remaining four loci (*rcc00555* and *rcc01866*) are frequently detected in bacterial genomes (Figs 1C and 2). Since within *Rhodobacterales* genes outside of ‘head–tail’ cluster evolve faster than their ‘head–tail’ counterparts (Hynes et al. 2016), we conjecture that many α -proteobacterial homologs of the former genes were not detected in our searches. Therefore, for the subsequent investigations we focused mostly on the analyses of the ‘head–tail’ locus genes.

3.3 Two distinct classes of RcGTA ‘head–tail’ homologs are present in α -proteobacteria

Within 187 α -proteobacterial genomes with at least one ‘head–tail’ cluster homolog, the RcGTA-like genes are dispersed across 474 genomic regions, 245 of which contain only 1–2 RcGTA homologs (Fig. 3). Due to significant similarity between RcGTA and genuine viral genes, some of these regions may belong to prophages unrelated to RcGTA. Indeed, 261 out of the 474 regions (55%) overlap with putative prophages identified using PhiSpy (Akhter et al. 2012) (Supplementary Fig. S3). However, PhiSpy also classified the RcGTA as a prophage, suggesting that

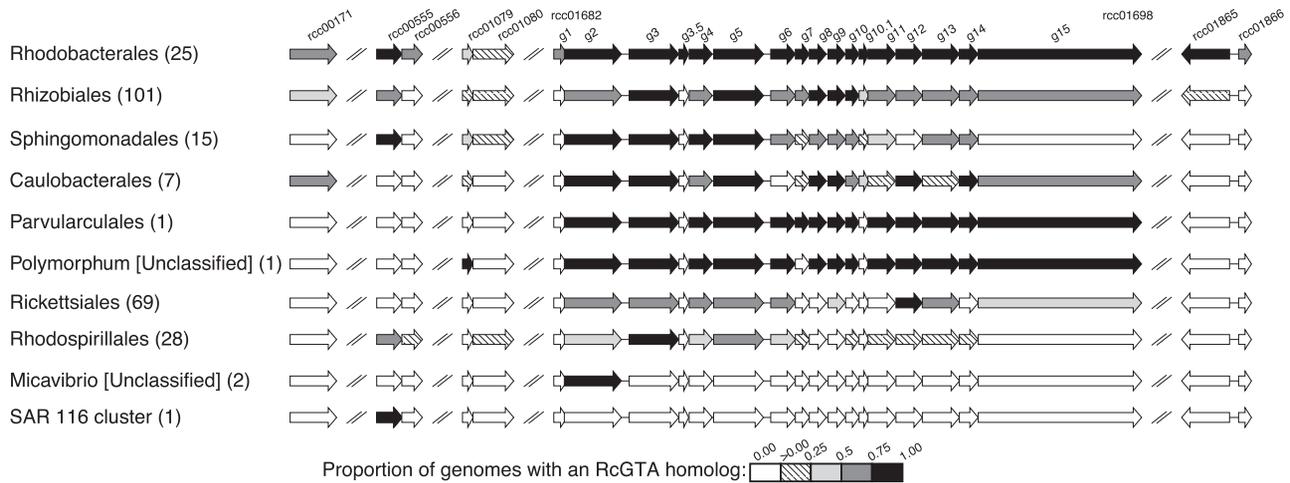


Figure 2. Presence of RcGTA-like genes in α -proteobacterial genomes. Each line represents the distribution of genes in the RcGTA genome within a taxonomic group. The number of surveyed genomes in a taxonomic group is listed in parentheses. The shades of gray represent the proportion of these genomes in which the RcGTA homologs are found (see figure inset). Genomic regions were visualized using R package genoplR (Guy et al. 2010).

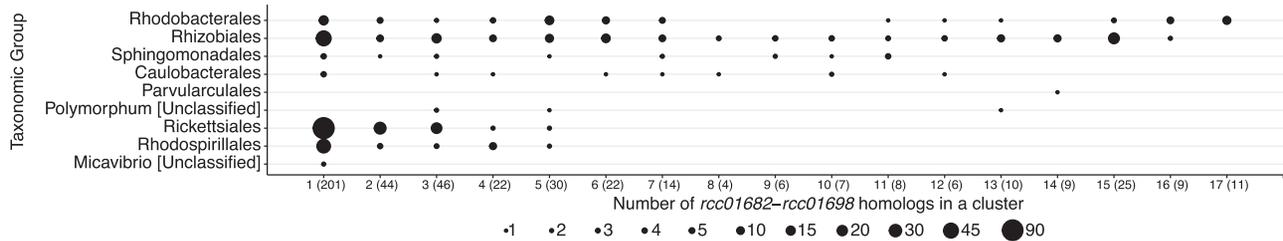


Figure 3. Size variation of the RcGTA-like 'head-tail' clusters among α -proteobacterial genomes. The number of homologs per genome (one to seventeen) is shown on the X-axis. The total number of genomes in each size category is shown in parentheses on the X-axis labels. The genomes are binned into taxonomic groups (arranged on Y-axis). The diameter of each circle is scaled with respect to the number of genomes (see inset for scale). Only 91 of the 474 clusters (19%) have at least nine RcGTA homologs. Sizes of RcGTA-like regions across several other bacterial taxonomic groups are shown in [Supplementary Fig. S2](#).

the predictions may include other false positives. Of 889 PhiSpy-predicted prophage regions within 255 α -proteobacterial genomes, only 535 (60%) are associated with an integrase gene—a gene expected to be present in a functional prophage but not in a GTA. RcGTA-like regions are overrepresented among the predicted prophages without an integrase gene (234 out of 354 [66%] vs. 27 out of 535 [5%]), hinting that PhiSpy might not be able to distinguish GTAs from genuine prophages. Also, PhiSpy did not detect small RcGTA-like regions due to the program's requirement of at least five consecutive genes with 'viral' functional annotations in a window of thirty genes. Therefore, we examined the RcGTA-like regions, and genes within them, for additional characteristics associated with viral genes—skewed nucleotide composition (Rocha and Danchin 2002) and faster substitution rates (Paterson et al. 2010; Drake 1999)—as well as the presence of neighboring viral genes.

Genes of viral origin generally have lower fraction of Gs and Cs (GC content) than their host genomes (Rocha and Danchin 2002). However, we do not observe this trend for the RcGTA-like regions. GC content of the regions with less than nine RcGTA homologs (hereafter referred as SCs) is, on average, equivalent to the GC content in the host (median of percent of relative change = -0.1% ; pairwise Wilcoxon test; P-value = 0.60), albeit with substantial variability across examined genomes (Fig. 4). On the other hand, GC content of the regions with at least nine RcGTA homologs (hereafter referred as LCs) is consistently, and in most cases significantly, higher than the GC content of their

host (median of percent of relative change = 6.57% ; pairwise t test; P-value <0.001) (Fig. 4). The elevated GC content of LCs also persists when compared to randomly sampled regions of the host genome of equivalent size (pairwise t-test; P-value <0.001). These observations prompted us to examine the features of LCs and SCs separately.

Using the PPD as a proxy for substitution rates, we found that LC genes evolve significantly slower than both their SC (all of the thirteen possible pairwise comparisons; Wilcoxon tests; P-values <0.0001) and viral homologs (fourteen out of the fifteen possible comparisons; Wilcoxon tests; P-values <0.005) (Supplementary Fig. S4 and Table S5). While many SC genes also evolve significantly slower than their viral homologs (six out of the ten possible comparisons; Wilcoxon tests; P-values <0.005), the 25–75 percentile range for PPDs of SC genes often overlapped with that of the viral homologs (Supplementary Fig. S4). Hence, SCs are more virus-like than LCs. This conjecture is supported by two additional observations. First, SCs are more frequently associated with unrelated viral genes than LCs (76 vs. 56%; Supplementary Fig. S3). In particular, SCs are more likely to reside in a vicinity of an integrase gene than LCs (45 out of 383 SCs vs. 1 out of 91 LCs). In some cases, SCs likely belong to viral elements that are clearly not GTAs: for example, in all analyzed genomes of the Rickettsiales' genera *Anaplasma* and *Ehrlichia* a singleton homolog of *rcc001695* (*g12*) is found within 1 kb of *virB* homologs, which in *Anaplasma phagocytophilum* belong to a virus-derived-type IV secretion system (Al-Khedery et al. 2012). Second, SCs

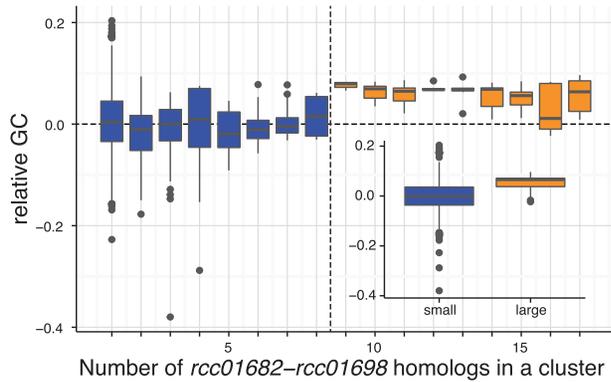


Figure 4. Difference in GC content between clusters of RCGTA ‘head-tail’ gene homologs and their host chromosome. The relative difference in GC content is calculated as $(GC_{\text{CGTA}} - GC_{\text{host}})/GC_{\text{host}}$ (Y-axis), and the data are presented by the number of RCGTA ‘head-tail’ homologs found in a cluster (X-axis). In the box-and-whisker plots the horizontal line marks the median; the boxes extend from the 25th to 75th percentile; the whiskers encompass values at most 1.5 times away from the 25th and 75th percentile range, and any value outside of the range is shown as a dot. The vertical dashed line separates two groups of clusters: ≤ 8 genes, or small clusters (SCs, in blue) and ≥ 9 genes, or large clusters (LC, in orange). The inset depicts the combined data from all small and LCs.

have very little conservation of flanking genes (Supplementary Fig. S5), while in LCs flanking genes are conserved within orders and sometimes even across larger phylogenetic distances (Supplementary Fig. S6).

Taken together, we hypothesize that large and SCs represent genomic regions under different selective pressures in the host chromosome and may have separate evolutionary histories. We further hypothesize that most SCs are viral elements unrelated to RCGTA and only LCs represent RCGTA-like elements.

3.4 RCGTA-like element was likely absent in the last common ancestor of α -proteobacteria

In what lineage did the RCGTA-like element originate? How did it propagate across α -proteobacteria? To address these questions, we examined the phylogenetic history of LCs. Within α -proteobacteria, LCs are found only within a monophyletic clade defined by Node 1 in Fig. 5 (hereafter referred as Clade 1), which consists of all analyzed representatives from the orders *Rhizobiales*, *Caulobacteriales*, *Rhodobacteriales*, *Parvularculales*, and *Sphingomonadales*. Two deeper branching α -proteobacterial clades, which include *Rickettsiales*, *Rhodospirillales*, and yet unclassified *Micavibrio* spp. and SAR 116 cluster bacteria, contain only SCs (Fig. 5). The largest of these SCs is made up of five genes, but most are singletons (Figs 2 and 3). RCGTA homologs are completely absent from *Pelagibacteriales* (formerly SAR11) and the basal α -proteobacterial order *Magnetococcales* (Fig. 5), although only five genomes from these two groups were available for the analyses.

It is possible that SCs represent decaying remnants of the previously functional GTAs. Bacterial genomes do not generally contain many pseudogenes due to genome streamlining (Lawrence et al. 2001). However, if pieces of RCGTA homologs that no longer form an ORF remain within the regions, they may still be detectable by similarity searches. Indeed, putative pseudogenes of RCGTA-like genes were detected in both large and SCs, but only 35% of them reside in SCs (Supplementary Figs S2, S5, and S6). Furthermore, only 3 of 199 SCs found in taxa outside of Clade 1 contain detectable putative pseudogenes of

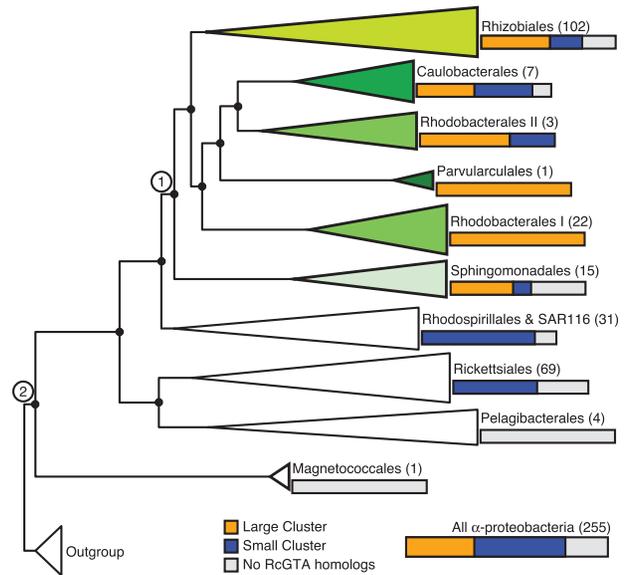


Figure 5. Distribution of RCGTA homologs in α -proteobacteria. The phylogenetic tree is a schematic version of the α -proteobacterial reference phylogeny, in which branches are collapsed at the taxonomic order level. The number of taxa in each collapsed clade is shown in parentheses. The tree is rooted using taxa from β -, γ -, and δ -proteobacteria (outgroup). The non-collapsed version of this tree is provided in Supplementary Fig. S3. The nodes with at least 60% bootstrap support are labeled with a black circle. Bars associated with each clade show the relative proportion of genomes that contain at least one LC (orange), only SC(s) (blue), or no detectable RCGTA homologs (gray). Nodes 1 and 2 define Clades 1 and 2 discussed in the text. Twenty-three of 150 genomes within Clade 1 and 32 out of the 105 genomes outside of Clade 1 have no detectable RCGTA homologs. While relationships among the α -proteobacterial orders remain debated (Lee et al. 2005; Gupta and Mok 2007; Williams et al. 2007; Thrash et al. 2011; Ferla et al. 2013; Luo 2014), the alternative histories do not affect the inferences of this study (data not shown). Scale bar, substitutions per site.

the RCGTA ‘head-tail’ genes (Supplementary Fig. S3), suggesting that SCs outside of Clade 1 are probably unrelated to RCGTA.

Taking into account the observation that the fraction of genomes with no detectable RCGTA homologs is higher in taxa outside than inside of Clade 1—30% versus 15% of the genomes, respectively—we hypothesize that the RCGTA-like element originate at the earliest on the branch leading to the last common ancestor of Clade 1 and not in the last common ancestor of α -proteobacteria (Clade 2 in Fig. 5), as was previously suggested (Lang et al. 2002; Lang and Beatty 2007).

3.5 Evolution of the RCGTA-like genome: vertically inherited or horizontally exchanged?

Did the RCGTA-like element appear on the branch leading to the last common ancestor of Clade 1, or did it spread across Clade 1 via HGT? Lang and Beatty (2007) inferred that RCGTA-like element had evolved vertically within α -proteobacteria, although only few genomes were available at the time. Within *Rhodobacteriales* RCGTA-like ‘head-tail’ clusters are also predicted to be vertically inherited, but phylogenetic trees of many individual genes were poorly resolved (Hynes et al. 2016). In contrast, homologs of *rcc00171* were likely horizontally exchanged within *Rhodobacteriales* (Hynes et al. 2016). To investigate the extent of HGT in the evolution of the RCGTA-like ‘head-tail’ genes across Clade 1, we reconstructed the phylogeny of LC genes and compared it to the reference tree of α -proteobacteria represented by a set of conserved α -proteobacterial genes.

Of eighty-three internal branches of the reference phylogeny, forty-two are in disagreement with the branching order of the LC-locus phylogeny (Fig. 6). Thirty-one of the forty-two conflicts (74%) are within taxonomic orders, with fourteen out of the twenty-one strongly supported conflicts limited to genera (Fig. 6 and Supplementary Fig. S7). Among the nine conflicts at deep branches, eight are due to an alternative position of *Methylobacterium nodulans* ORS 2060 (Fig. 6 and Supplementary Fig. S7), which, in addition to grouping outside of its order Rhizobiales, has five non-identical LCs (Supplementary Fig. S6). Therefore, evolution of ‘head–tail’ locus is likely impacted by HGT events, but most of them have occurred between closely related taxa.

Intriguingly, at least occasionally LCs have also been transferred across very large evolutionary distances. While LCs are predominantly found in α -proteobacteria, the genomes of the actinobacteria *Streptomyces purpurogeneiscleroticus* NRRL B-2952 and *Asanoa ferruginea* NRRL B-16430, the cyanobacterium *Scytonema millei* VB511283, and the γ -proteobacterium *Pseudomonas bauzanensis* W13Z2 also contain one LC each (Supplementary Figs S2 and S6). These four non- α -proteobacterial taxa branch within the Clade 1 in three separate locations (Fig. 6), indicating that these taxa likely acquired the LCs from the α -proteobacteria at least three times independently. Moreover, genes immediately upstream and downstream of *S. millei* and *P. bauzanensis*’ LCs are homologous to the genes that flank LCs in their respective α -proteobacterial sister taxa (Supplementary Fig. S6). These flanking genes are also of α -proteobacterial origin (Supplementary Fig. S8), and hence were likely acquired in the same HGT events.

The duration of RcGTA-like elements’ association with bacterial genomes can also be gauged from the stability of the host chromosome regions that surround LCs. For example, highly mobile genes, such as transposable elements and some prophages, are often found in genomic islands (Juhás et al. 2009)—dynamic genomic regions that exhibit poor gene synteny even across short phylogenetic distances (Rodríguez-Valera et al. 2016). On the other hand, gene synteny of stable regions is preserved across larger evolutionary distances, although it is still subject to decay due to gene rearrangements (Rocha 2006; Brilli et al. 2013). LCs found in the same taxonomic order are generally flanked by the same, conserved α -proteobacterial genes (Supplementary Fig. S6), hinting that they might be part of stable regions of bacterial chromosomes and not within genomic islands. We modeled gene order decay upstream and downstream of LC loci as a function of phylogenetic distance (Rocha 2006) and compared the conservation of genes flanking LCs with that of a transposable element from IS911/IS3 family and two operons of conserved housekeeping genes, the ribosomal protein operon and the ATP synthase operon. The IS911/IS3 element can undergo non-targeted integration (Chandler et al. 2015), and thus, genes adjoining these transposons are not expected to be homologous even across closely related taxa. In contrast, gene order near presumably non-mobile housekeeping operons are expected to be much more conserved across time. Indeed, we found that the regions surrounding the transposable element are too variable to even fit the model, while the decay of gene synteny near LCs is similar to that of the two operons (Supplementary Fig. S9). Therefore, it is unlikely that RcGTA-like elements reside in dynamic regions of their host genomes.

Although sequence divergence of the remaining four loci of the RcGTA genome makes it difficult to track their evolutionary history, in sixty-six α -proteobacterial genomes with detected homologs of *rcc00171* and *rcc00555*, twenty-two LCs from

Rhizobiales and Rhodobacterales have *rcc00171* homologs at their 3’ end and six of the twenty-two LCs additionally have *rcc00555* immediately downstream of an *rcc00171* homolog (Supplementary Fig. S6). This anecdotal evidence suggests that in the past RcGTA-like elements may have carried at least some of the genes in one locus with the ‘head–tail’ genes, and that the division of the GTA genome into multiple loci may have happened after the bacterium-GTA association was already established.

Taken together, our analyses suggest a complex evolutionary history of the RcGTA-like elements. The congruence of the reference and the LC-locus phylogenies at taxonomic order and family levels and conservation of LC-flanking genes within orders suggest that RcGTA-like element was likely acquired before the contemporary orders within Clade 1 have diversified, and since that time the element co-existed with the clade members. However, the congruence does not automatically equate to vertical inheritance, since even genes highly conserved across α -proteobacteria are not immune to HGT and thus the reference phylogeny may not represent the strictly vertical history of α -proteobacteria (Gogarten et al. 2002; Andam et al. 2010). Indeed, strong phylogenetic conflicts between the reference and LC locus phylogenies of Clade 1 at both deep and shallow branches, and presence of LCs in several taxa outside of the α -proteobacteria indicate that HGT has been shaping the evolution of RcGTA-like element to a non-negligible degree. Since the majority of the strongly supported conflicts are found within genera, HGT probably played a larger role in the recent evolution of RcGTA-like elements than in their dissemination across Clade 1.

4. Discussion

Our survey of RcGTA homologs in the genomes of bacteria and bacterial viruses provide new insights into fascinatingly intertwined evolutionary histories of virus-like elements and their bacterial hosts. Although we infer that the RcGTA-like element is not as ancient as the class α -proteobacteria itself, RcGTA homologs are ubiquitous across a clade that spans multiple α -proteobacterial orders (Clade 1 in Fig. 5) and, according to molecular clock estimates, has originated between 777 Mya (Luo et al. 2013) and 1710 Mya (with a 95% CI of 1977–1402 Ma; Gregory Fournier, pers. comm.). Therefore, either RcGTA-like element has been associated with a group of α -proteobacteria for hundreds of millions of years, or it is a vestige of a very successful virus capable of inter-species and inter-order infections. The latter scenario is supported by experimentally observed long-range HGT via GTAs produced by *Roseovirus nubinhibens* ISM (Mcdaniel et al. 2010). However, lack of CRISPR repeats that match RcGTA genome, absence of closely related contemporary viruses in GenBank, and order-level congruence of the evolutionary histories of RcGTA and conserved α -proteobacterial genes point against the GTA being derived from a broad host-range virus that recently invaded this clade. Instead, our data are more compatible with a long-term association between a GTA and its bacterial host.

Over this time, the GTA genome has gone through extensive changes, as hinted by the putative HGT events within taxonomic orders, presence of likely pseudogenized RcGTA homologs within a few putative GTAs, and division of the RcGTA genome into multiple loci. Could these modifications simply mean that this RcGTA-like element is just a defective prophage (Solió and Marrs 1977; Redfield 2001) not yet purged from the genomes? Since bacterial genomes are generally under deletion bias of unnecessary DNA (Mira et al. 2001), conservation of

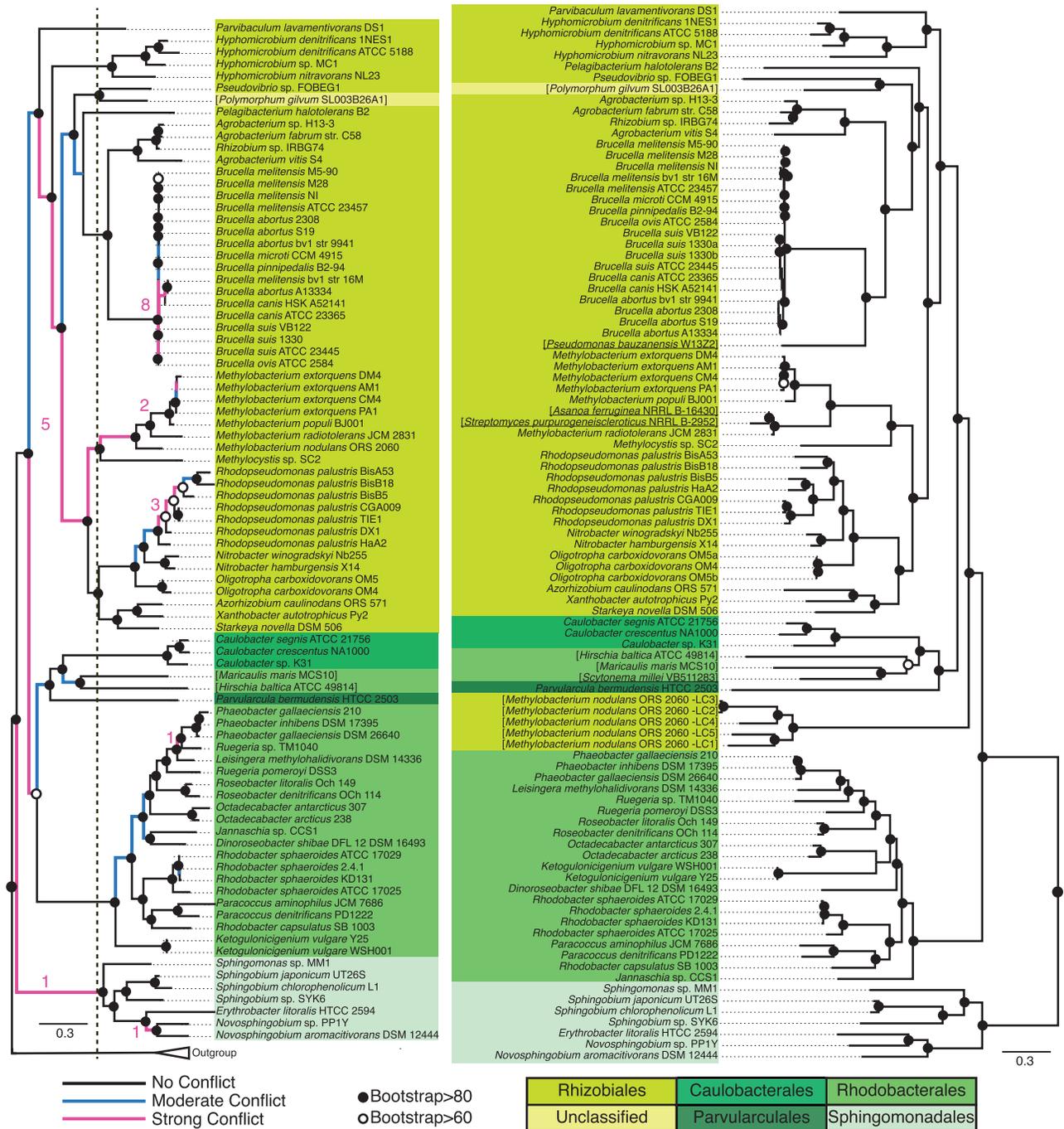


Figure 6. Comparison of the phylogenetic history of LCs (right) to the reference α -proteobacterial phylogeny (left). The branches of the reference tree are colored according to the strength of conflict, as measured by an Internode Certainty (IC) value (Salichos et al. 2014). IC values for all branches are shown in Supplementary Fig. S7. The overall number of branches with strongly conflicting IC values is shown as numbers (in pink) for 'deep' and 'shallow' branches, as demarcated by a vertical dashed line. Fourteen of the fifteen strong conflicts at shallow branches occur within genera, and five of the six strong conflicts at deep branches are due to the different position of *Methylobacterium nodulans* ORS 2060. The phylogenetic tree on the left is rooted with the same outgroup as in Supplementary Fig. S3, while the phylogenetic tree on the right should be considered unrooted. Underlined taxa names represent four non- α -proteobacterial genomes that contain an LC (see also Supplementary Figs S6 and S9). Taxa names in brackets represent lineages that branch within a taxonomic order different from the one designated by the NCBI Taxonomy database (NCBI Resource Coordinators 2017). The tree heights (i.e. the sum of all branch lengths of a phylogenetic tree) of the LC-locus and reference phylogenies are 31.19 and 19.39 substitutions per site, respectively, suggesting that LC locus evolves ~ 1.6 times faster than the conserved α -proteobacterial genes. Scale bar, substitutions per site.

'junk' regions across over millions of years is unexpected. Moreover, given the universal mutation bias towards AT-rich DNA (Hershberg and Petrov 2010), over such prolonged time we would expect the defective prophage region to become more

AT-rich than the host DNA, and to accumulate mutations evenly across non-synonymous and synonymous sites of ORFs (Andersson and Andersson 2001). Yet, genomic regions that encode LCs are not AT-richer than the rest of the host

chromosomes, and, at least across *Rhodobacterales*, RcGTA homologs do not exhibit increased rates of non-synonymous substitutions expected of pseudogenes (Lang et al. 2012), suggesting that LCs are unlikely to be decaying prophages.

Surprisingly, RcGTA-like gene clusters have a significantly elevated % G + C than chromosomes of their hosts. Evolution of GC content and its variation across and within the prokaryotic genomes is an unsolved puzzle, with many possibilities invoked to explain it (Hildebrand et al. 2010; Rocha and Feil 2010; Lassalle et al. 2015). In one hypothesis, AT-richness of highly expressed genes is explained by accumulation of additional mutations during the single-stranded DNA state required for transcription (Hildebrand et al. 2010). Since RcGTA-producing cells lyse and thus leave no progeny, the GTA genes are never expressed by cells that reproduce. Hence, it is tempting to speculate that selection for reduced gene expression may have played a role in driving the GC content of the ‘head-tail’ locus higher than that of the expressed host genes. Future testing of this hypothesis and exploring other possible explanations are necessary to account for the aberrant GC content of the RcGTA-like regions within Clade 1 genomes.

But if RcGTA-like elements are indeed functional, why do we observe putatively pseudogenized and apparently incomplete gene repertoires of many LCs when compared with the RcGTA genome? Perhaps, genomes of these elements are also split into different pieces, have some components replaced from sources unrelated to the RcGTA, or have evolved to a different functionality than RcGTA. For example, GTA gene expression and GTA particle release was demonstrated for *Ruegeria pomeroyi* DSS-3 (Biers et al. 2008) and *Rhodovulum sulfidophilum* (Nagao et al. 2015), suggesting that GTA is functional in these lineages. However, the *R. pomeroyi* and *R. sulfidophilum* genome contain no detectable homologs of *g1*, and this gene is essential for RcGTA production in *R. capsulatus* (Hynes et al. 2016). Therefore, even presumably ‘incomplete’ LCs may encode functional GTAs—a conjecture that awaits experimental demonstrations.

With the evidence pointing against the decaying or a selfish nature of the RcGTA-like elements, and given that cells that produce RcGTA lyse (Westbye et al. 2013) and, therefore, the selection for the maintenance of RcGTA cannot occur at a level of an organism, earlier suggestions that GTA maintenance is due to some population-level benefits (Lang et al. 2012) remain most plausible explanation of persistence of RcGTA-like elements in the Clade 1 genomes. However, specific advantage(s) associated with acquisition of ~4 kb pieces of random DNA via GTA remain to be elucidated. Among the proposed drivers of GTA maintenance are effective exchange of useful genes among the population of heterogeneous cells (Smillie et al. 2011), improved response to environmental stressors, such as DNA damage (Marrs et al. 1977; Brimacombe et al. 2015) and nutritional deficits (Lang et al. 2012). Possibilities unrelated to GTA-mediated DNA delivery, such as protection of the host population against infection by other viruses via the superinfection immunity mechanism (Díaz-Muñoz 2017), and increase of the host population resilience against various environmental stressors such as antibiotics (Wang et al. 2010), should also be considered. Perhaps, there is no single benefit associated with RcGTA-like elements across all Clade 1 taxa. Instead, different lineages may experience selection pressures unique to their ecological niches.

Interestingly, GTA-associated benefits have been proposed to drive evolutionary innovations of cellular lifeforms that go beyond population level. For example, it has been hypothesized that abundance and mixed ancestry of α -proteobacterial genes

in eukaryotic genomes is due to the delivery of such genes by RcGTA-like elements into proto-eukaryote genome, and that such ‘seeding’ facilitated the later integration of mitochondrial progenitor and the proto-eukaryotic host (Richards and Archibald 2011). Our study, however, does not support this hypothesis. First, based on the estimated age of the Clade 1 (777–1710 Mya) and the appearance of eukaryotic, and hence mitochondria-containing, organisms in the fossil record (1600 and possibly 1800 Mya; Knoll 2014), the RcGTA-like system likely originated after the eukaryogenesis took place. Second, α -proteobacterial lineage(s) suggested to have given rise to mitochondrial progenitor (Brindefalk et al. 2011; Thrash et al. 2011; Wang and Wu 2015) appear to lack LCs and therefore are unlikely to have had RcGTA-like elements. In another intriguing hypothesis, an unrelated GTA in *Bartonella* has been proposed to have facilitated adaptive radiation within the genus (Guy et al. 2013). Perhaps, acquisition of the RcGTA-like element by the last common ancestor of Clade 1 taxa represent a similar innovation that led to diversification of this clade and success of its members in a variety of ecological settings that span soil, freshwater, marine, waste water, wetland, and eukaryotic intracellular and extracellular environments.

Despite the appearance that bacteria maintain GTA genes for the potential contribution to their own evolutionary success, there is an undeniable shared evolutionary history between RcGTA-like elements and *bona fide* viruses. Presence of homologs for most RcGTA-like genes in viral genomes, as well as similarity in the organization of the ‘head-tail’ cluster and corresponding region of a typical siphovirus genome (Huang et al. 2011), led to a prevailing hypothesis that RcGTA originated in an initial co-option of a virus and evolved in subsequent modification of the progenitor prophage genome via vertical descent and HGT (Lang et al. 2017). Our data are compatible with this scenario, and further suggest an even larger role of HGT in the evolution of RcGTA-like elements than currently acknowledged, shaping GTAs both within α -proteobacterial orders, but also at least occasionally transferring RcGTA-like regions across the higher-order α -proteobacterial taxa and even to different phyla. Yet, there are hints of even more entangled connection between RcGTA-like genes found in cellular and viral genomes: a few viruses that infect *Paracoccus*, *Roseobacter*, *Caulobacter*, and *Rhodobacter* spp., and a yet unclassified marine *Rhizobiales* str. JL001 bacterium, carry RcGTA-like genes that are phylogenetically placed within their cellular homologs (Zhan et al. 2016), suggesting that some viruses acquired genes from the cellular RcGTA-like regions. Given recent propositions of (1) cellular origin of even such ‘hallmark’ viral genes like those encoding capsid proteins (Krupovic and Koonin 2017), (2) the possibly virus-independent origin of cellular microcompartments that resemble viral structures (Bobik et al. 2015), and (3) the potentially primordial origin of type VI secretion system that resembles a bacteriophage tail (Böck et al. 2017), we should not exclude a possibility that genomic regions encoding RcGTA-like nanostructures did not originate by a co-option of a prophage and its subsequent modification, but instead are a mosaic originally assembled within bacterial genomes from the available cellular and viral parts.

Authors’ contributions

M.S. and O.Z. designed the study. M.S. and S.M.S. performed the analyses. M.S., S.M.S., and O.Z. wrote the manuscript. All authors have seen and approved the final manuscript.

Data availability

Amino acid sequences of RcGTA homologs in bacterial and viral genomes; amino acid sequence alignments of RcGTA homologs from the LCs, SCs, and viruses; concatenated alignment of ninety-nine conserved α -proteobacterial genes and RcGTA homologs from LCs; and discussed phylogenetic trees in Newick format are available via FigShare at DOI 10.6084/m9.figshare.5406733.

Supplementary data

Supplementary data are available at *Virus Evolution* online.

Acknowledgements

We would like to thank Andrew S. Lang, J. Thomas Beatty, and Rosemary J. Redfield for numerous stimulating discussions regarding GTA evolution.

Funding

This work was supported by the National Science Foundation (NSF-DEB 1551674 to O.Z.); the Simons Foundation (Investigator in Mathematical Modeling of Living Systems award 327936 to O.Z.); the Neukom Institute CompX award to O.Z.; and Dartmouth start-up funds to O.Z.

Conflict of interest: None declared.

References

- Aberer, A. J., Krompass, D., and Stamatakis, A. (2013) 'Pruning rogue taxa improves phylogenetic accuracy: an efficient algorithm and web-service', *Systematic Biology*, 62: 162–6
- Akhter, S., Aziz, R. K., and Edwards, R. A. (2012) 'PhiSpy: a novel algorithm for finding prophages in bacterial genomes that combines similarity- and composition-based strategies', *Nucleic Acids Research*, 40: e126
- Al-Khedery, B. et al. (2012) 'Structure of the type IV secretion system in different strains of *Anaplasma phagocytophilum*', *BMC Genomics*, 13: 678
- Alkhnabashi, O. S. et al. (2016) 'Characterizing leader sequences of CRISPR loci', *Bioinformatics*, 32: i576–85
- Altschul, S. F. et al. (1997) 'Gapped BLAST and PSI-BLAST: a new generation of protein database search programs', *Nucleic Acids Research*, 25: 3389–402
- Andam, C. P., Williams, D., and Gogarten, J. P. (2010) 'Biased gene transfer mimics patterns created through shared ancestry', *Proceedings of the National Academy of Sciences of the United States of America*, 107: 10679–84
- Andersson, J. O., and Andersson, S. G. (2001) 'Pseudogenes, junk DNA, and the dynamics of Rickettsia genomes', *Molecular Biology and Evolution*, 18: 829–39.
- Bates, D. M., and Watts, D. G. (1988). *Nonlinear Regression Analysis and its Applications*. New York: Wiley & Sons, Inc.
- Berger, S. A., Krompass, D., and Stamatakis, A. (2011) 'Performance, accuracy, and Web server for evolutionary placement of short sequence reads under maximum likelihood', *Systematic Biology*, 60: 291–302
- Berglund, E. et al. (2009) 'Run-off replication of host-adaptability genes is associated with gene transfer agents in the genome of mouse-infecting *Bartonella grahamii*', *PLoS Genetics*, 5: e1000546
- Bertani, G. (1999) 'Transduction-like gene transfer in the methanogen *Methanococcus voltae*', *Journal of Bacteriology*, 181: 2992–3002
- Biers, E. J. et al. (2008) 'Occurrence and expression of gene transfer agent genes in marine bacterioplankton', *Applied and Environmental Microbiology*, 74: 2933–9
- Bobay, L. M., Touchon, M., and Rocha, E. P. (2014) 'Pervasive domestication of defective prophages by bacteria', *Proceedings of the National Academy of Sciences of the United States of America*, 111: 12127–32
- Bobik, T. A., Lehman, B. P., and Yeates, T. O. (2015) 'Bacterial microcompartments: widespread prokaryotic organelles for isolation and optimization of metabolic pathways', *Molecular Microbiology*, 98: 193–207
- Böck, D. et al. (2017) 'In situ architecture, function, and evolution of a contractile injection system', *Science*, 357: 713–7
- Borgeaud, S. et al. (2015) 'Bacterial evolution. The type VI secretion system of *Vibrio cholerae* fosters horizontal gene transfer', *Science*, 347: 63–7.
- Brilli, M. et al. (2013) 'Short and long-term genome stability analysis of prokaryotic genomes', *BMC Genomics*, 14: 309
- Brimacombe, C. A. et al. (2015) 'Homologues of genetic transformation DNA import genes are required for *Rhodobacter capsulatus* gene transfer agent recipient capability regulated by the response regulator CtrA', *Journal of Bacteriology*, 197: 2653–63.
- Brimacombe, C. et al. (2013) 'Quorum-sensing regulation of a capsular polysaccharide receptor for the *Rhodobacter capsulatus* gene transfer agent (RcGTA)', *Molecular Microbiology*, 87: 802–17.
- Brindefalk, B. et al. (2011) 'A phylometagenomic exploration of oceanic alphaproteobacteria reveals mitochondrial relatives unrelated to the SAR11 clade', *PLoS One*, 6: e24457
- Canchaya, C. et al. (2003) 'Prophage genomics', *Microbiology and Molecular Biology Reviews*, 67: 238–76
- Chandler, M. et al. (2015) 'Copy-out-Paste-in transposition of IS911: A major transposition pathway', *Microbiology Spectrum*, 3, doi:10.1128/microbiolspec.MDNA3-0031-2014.
- Chen, J. C., and Stephens, C. (2007) 'Bacterial cell cycle: completing the circuit', *Current Biology*, 17: R203–6
- Cobián Güemes, A. G. et al. (2016) 'Viruses as winners in the game of life', *Annual Review of Virology*, 3: 197–214
- Desantis, T. et al. (2006) 'Greengenes, a chimera-checked 16S rRNA gene database and workbench compatible with ARB', *Applied and Environmental Microbiology*, 72: 5069–72
- Díaz-Muñoz, S. L. (2017) 'Viral coinfection is shaped by host ecology and virus-virus interactions across diverse microbial taxa and environments', *Virus Evolution*, 3: vex011
- van Dongen, S. M. (2000). 'Graph Clustering by Flow Simulation', PhD thesis, University of Utrecht.
- Drake, J. W. (1999) 'The distribution of rates of spontaneous mutation over viruses, prokaryotes, and eukaryotes', *Annals of the New York Academy of Sciences*, 870: 100–7
- Eddy, S. R. (1998) 'Profile hidden Markov models', *Bioinformatics*, 14: 755–63
- Edgar, R. (2004) 'MUSCLE: a multiple sequence alignment method with reduced time and space complexity', *BMC Bioinformatics*, 5: 113
- Edgar, R. (2007) 'PILER-CR: fast and accurate identification of CRISPR repeats', *BMC Bioinformatics*, 8: 18
- Ferla, M. et al. (2013) 'New rRNA gene-based phylogenies of the alphaproteobacteria provide perspective on major groups, mitochondrial ancestry and phylogenetic instability', *PLoS One*, 8: e83383
- Fogg, P. C., Westbye, A., and Beatty, J. T. (2012) 'One for all or all for one: heterogeneous expression and host cell lysis are key

- to gene transfer agent activity in *Rhodobacter capsulatus*', *PLoS One*, 7: e43772.
- Frost, L. et al. (2005) 'Mobile genetic elements: the agents of open source evolution', *Nature Reviews of Microbiology*, 3: 722–32
- Gogarten, J. P., Doolittle, W. F., and Lawrence, J. G. (2002) 'Prokaryotic evolution in light of gene transfer', *Molecular Biology and Evolution*, 19: 2226–38
- Grazziotin, A. L., Koonin, E. V., and Kristensen, D. M. (2017) 'Prokaryotic virus orthologous groups (pVOGs): a resource for comparative genomics and protein family annotation', *Nucleic Acids Research*, 45: D491–8
- Gupta, R. S., and Mok, A. (2007) 'Phylogenomics and signature proteins for the alpha proteobacteria and its main groups', *BMC Microbiology*, 7: 106
- Guy, L., Kultima, J. R., and Andersson, S. G. E. (2010) 'genoPlotR: comparative gene and genome visualization in R', *Bioinformatics*, 26: 2334–5
- , ——, and —— (2013) 'A gene transfer agent and a dynamic repertoire of secretion systems hold the keys to the explosive radiation of the emerging pathogen *Bartonella*', *PLoS Genetics*, 9: e1003393.
- Hershberg, R., and Petrov, D. A. (2010) 'Evidence that mutation is universally biased towards AT in bacteria', *PLoS Genetics*, 6: e1001115
- He, Z. et al. (2016) 'Evolview v2: an online visualization and management tool for customized and annotated phylogenetic trees', *Nucleic Acids Research*, 44: W236–41
- Hildebrand, F., Meyer, A., and Eyre-Walker, A. (2010) 'Evidence of selection upon genomic GC-content in bacteria', *PLoS Genetics*, 6: e1001107
- Huang, S. et al. (2011) 'Complete genome sequence of a marine roseophage provides evidence into the evolution of gene transfer agents in alphaproteobacteria', *Virology Journal*, 8: 124
- Humphrey, S. B. et al. (1997) 'Purification and characterization of VSH-1, a generalized transducing bacteriophage of *Serpulina hydysenteriae*', *Journal of Bacteriology*, 179: 323–9.
- Hynes, A. P. et al. (2012) 'DNA packaging bias and differential expression of gene transfer agent genes within a population during production and release of the *Rhodobacter capsulatus* gene transfer agent, RcGTA', *Molecular Microbiology*, 85: 314–25.
- et al. (2016) 'Functional and evolutionary characterization of a gene transfer agent's multilocus "genome"', *Molecular Biology and Evolution*, 33: 2530–43
- Iranzo, J., Krupovic, M., and Koonin, E. V. (2016) 'The double-stranded DNA virosphere as a modular hierarchical network of gene sharing', *mBio*, 7: e00978-16
- & —— and —— (2017) 'A network perspective on the virus world', *Communicative & Integrative Biology*, 10: e1296614
- Juhas, M. et al. (2009) 'Genomic islands: tools of bacterial horizontal gene transfer and evolution', *FEMS Microbiology Reviews*, 33: 376–93
- Ju, K.-S. et al. (2015) 'Discovery of phosphonic acid natural products by mining the genomes of 10, 000 actinomycetes', *Proceedings of the National Academy of Sciences of the United States of America*, 112: 12175–80
- Knoll, A. H. (2014) 'Paleobiological perspectives on early eukaryotic evolution', *Cold Spring Harbor Perspectives in Biology*, 6: a016121
- Koonin, E., Senkevich, T., and Dolja, V. V. (2006) 'The ancient virus world and evolution of cells', *Biology Direct*, 1: 29
- Koonin, E. V. (2016) 'Horizontal gene transfer: essentiality and evolvability in prokaryotes, and roles in evolutionary transitions', *F1000Research*, 5: 1805
- Kristensen, D. M. et al. (2013) 'Orthologous gene clusters and taxon signature genes for viruses of prokaryotes', *Journal of Bacteriology*, 195: 941–50
- Krupovic, M., and Koonin, E. V. (2017) 'Multiple origins of viral capsid proteins from cellular ancestors', *Proceedings of the National Academy of Sciences of the United States of America*, 114: E2401–10
- Kuchinski, K. S. et al. (2016) 'The SOS response master regulator LexA regulates the gene transfer agent of *Rhodobacter capsulatus* and represses transcription of the signal transduction protein CckA', *Journal of Bacteriology*, 198: 1137–48.
- Lang, A. S., and Beatty, J. T. (2000) 'Genetic analysis of a bacterial genetic exchange element: the gene transfer agent of *Rhodobacter capsulatus*', *Proceedings of the National Academy of Sciences of the United States of America*, 97: 859–64
- , and —— (2007) 'Importance of widespread gene transfer agent genes in a-proteobacteria', *Trends in Microbiology*, 15: 54–62
- , Taylor, T., and Beatty, J. T. (2002) 'Evolutionary implications of phylogenetic analyses of the gene transfer agent (GTA) of *Rhodobacter capsulatus*', *Journal of Molecular Evolution*, 55: 534–43.
- , Westbye, A. B., and —— (2017) 'The distribution, evolution, and roles of gene transfer agents in prokaryotic genetic exchange', *Annual Review of Virology*, 4: 87–104
- , Zhaxybayeva, O., and —— (2012) 'Gene transfer agents: phage-like elements of genetic exchange', *Nature Reviews. Microbiology*, 10: 472–82
- Lassalle, F. et al. (2015) 'GC-content evolution in bacterial genomes: the biased gene conversion hypothesis expands', *PLoS Genetics*, 11: e1004941.
- Lawrence, J. G., Hendrix, R. W., and Casjens, S. (2001) 'Where are the pseudogenes in bacterial genomes?', *Trends in Microbiology*, 9: 535–40
- Lee, K. B. et al. (2005) 'The hierarchical system of the "Alphaproteobacteria": description of Hyphomonadaceae fam. nov., Xanthobacteraceae fam. nov. and Erythrobacteraceae fam. nov.', *International Journal of Systematic and Evolutionary Microbiology*, 55: 1907–19
- Le, S. Q., and Gascuel, O. (2008) 'An improved general amino acid replacement matrix', *Molecular Biology and Evolution*, 25: 1307–20
- Leung, M. et al. (2012) 'The GtaR protein negatively regulates transcription of the gtaRI operon and modulates gene transfer agent (RcGTA) expression in *Rhodobacter capsulatus*', *Molecular Microbiology*, 83: 759–74.
- et al. (2010) 'The gene transfer agent of *Rhodobacter capsulatus*', *Advances in Experimental Medicine and Biology*, 675: 253–64
- Li, L., Stoekert, C. J., Jr., and Roos, D. S. (2003) 'OrthoMCL: Identification of ortholog groups for eukaryotic genomes', *Genome Research*, 13: 2178–89.
- Luo, H. (2014) 'Evolutionary origin of a streamlined marine bacterioplankton lineage', *The ISME Journal*, 9: 1423–33
- et al. (2013) 'Evolution of divergent life history strategies in marine alphaproteobacteria', *mBio*, 4: e00373–13
- Makarova, K. S. et al. (2015) 'An updated evolutionary classification of CRISPR-Cas systems', *Nature Reviews Microbiology*, 13: 722–36
- Mann, T. H. et al. (2016) 'A cell cycle kinase with tandem sensory PAS domains integrates cell fate cues', *Nature Communications*, 7: 11454
- Marrs, B. (1974) 'Genetic recombination in *Rhodopseudomonas capsulata*', *Proceedings of the National Academy of Sciences of the United States of America*, 71: 971–3.

- , Wall, J. D., and Gest, H. (1977) 'Emergence of the biochemical genetics and molecular biology of photosynthetic bacteria', *Trends in Biochemical Sciences*, 2: 105–8
- Mcdaniel, L. et al. (2010) 'High frequency of horizontal gene transfer in the oceans', *Science*, 330: 50
- et al. (2012) 'Environmental factors influencing gene transfer agent (GTA) mediated transduction in the subtropical ocean', *PLoS One*, 7: e43506
- McHugh, C. A. et al. (2014) 'A virus capsid-like nanocompartment that stores iron and protects bacteria from oxidative stress', *The EMBO Journal*, 33: 1896–911
- Mercer, R. G., and Lang, A. S. (2014) 'Identification of a predicted partner-switching system that affects production of the gene transfer agent RcGTA and stationary phase viability in *Rhodobacter capsulatus*', *BMC Microbiology*, 14: 71
- Mira, A., Ochman, H., and Moran, N. A. (2001) 'Deletional bias and the evolution of bacterial genomes', *Trends in Genetics*, 17: 589–96
- Nagao, N. et al. (2015) 'The gene transfer agent-like particle of the marine phototrophic bacterium *Rhodovulum sulfidophilum*', *Biochemistry and Biophysics Reports*, 4: 369–74
- NCBI Resource Coordinators. (2017) 'Database resources of the National Center for Biotechnology Information', *Nucleic Acids Research*, 45: D12–7
- Paterson, S. et al. (2010) 'Antagonistic coevolution accelerates molecular evolution', *Nature*, 464: 275–8
- Price, M. N., Dehal, P. S., Arkin, A. P., and Poon, A. F. Y. (2010) 'FastTree 2—approximately maximum-likelihood trees for large alignments', *PLoS One*, 5: e9490
- Rapp, B. J., and Wall, J. D. (1987) 'Genetic transfer in *Desulfovibrio desulfuricans*', *Proceedings of the National Academy of Sciences of the United States of America*, 84: 9128–30.
- Redfield, R. J. (2001) 'Do bacteria have sex?' *Nature Reviews Genetics*, 2: 634–9
- Richards, T. A., and Archibald, J. M. (2011) 'Cell evolution: gene transfer agents and the origin of mitochondria', *Current Biology*, 21: R112–4
- Rocha, E. P. C. (2006) 'Inference and analysis of the relative stability of bacterial chromosomes', *Molecular Biology and Evolution*, 23: 513–22
- and Danchin, A. (2002) 'Base composition bias might result from competition for metabolic resources', *Trends in Genetics*, 18: 291–4
- and Feil, E. J. (2010) 'Mutational patterns cannot explain genome composition: Are there any neutral sites in the genomes of bacteria?' *PLoS Genetics*, 6: e1001104
- Rodríguez-Valera, F., Martín-Cuadrado, A.-B., and López-Pérez, M. (2016) 'Flexible genomic islands as drivers of genome evolution', *Current Opinion in Microbiology*, 31: 154–60
- Salichos, L., Stamatakis, A., and Rokas, A. (2014) 'Novel information theory-based measures for quantifying incongruence among phylogenetic trees', *Molecular Biology and Evolution*, 31: 1261–71
- Sarris, P. F. et al. (2014) 'A phage tail-derived element with wide distribution among both prokaryotic domains: a comparative genomic and phylogenetic study', *Genome Biology and Evolution*, 6: 1739–47
- Schaefer, A. L. et al. (2002) 'Long-chain acyl-homoserine lactone quorum-sensing regulation of *Rhodobacter capsulatus* gene transfer agent production', *Journal of Bacteriology*, 184: 6515–21.
- Schloss, P. et al. (2009) 'Introducing mothur: open-source, platform-independent, community-supported software for describing and comparing microbial communities', *Applied and Environmental Microbiology*, 75: 7537–41
- Selengut, J. D. et al. (2007) 'TIGRFAMs and Genome Properties: tools for the assignment of molecular function and biological process in prokaryotic genomes', *Nucleic Acids Research*, 35/Database: D260–4
- Sen, D. et al. (2015) 'Draft genome sequence of the terrestrial cyanobacterium *Scytonema millei* VB511283, isolated from Eastern India', *Genome Announcements*, 3: e00009–15
- Shikuma, N. J. et al. (2014) 'Marine tubeworm metamorphosis induced by arrays of bacterial phage tail-like structures', *Science*, 343: 529–33
- Smillie, C. S. et al. (2011) 'Ecology drives a global network of gene exchange connecting the human microbiome', *Nature*, 480: 241–4
- Soliz, M., and Marrs, B. (1977) 'The gene transfer agent of *Rhodospseudomonas capsulata*. Purification and characterization of its nucleic acid', *Archives of Biochemistry and Biophysics*, 181: 300–7.
- Stamatakis, A. (2014) 'RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies', *Bioinformatics*, 30: 1312–3
- Sutter, M. et al. (2008) 'Structural basis of enzyme encapsulation into a bacterial nanocompartment', *Nature Structural & Molecular Biology*, 15: 939–47
- Thrash, J. et al. (2011) 'Phylogenomic evidence for a common ancestor of mitochondria and the SAR11 clade', *Scientific Reports*, 1: 1–9
- Touchon, M., Bernheim, A., and Rocha, E. P. (2016) 'Genetic and life-history traits associated with the distribution of prophages in bacteria', *The ISME Journal*, 10: 2744–54
- , Moura de Sousa, J. A., and —— (2017) 'Embracing the enemy: the diversification of microbial gene repertoires by phage-mediated horizontal gene transfer', *Current Opinion in Microbiology*, 38: 66–73
- Wang, X. et al. (2014) 'Draft genome sequence of halotolerant polycyclic aromatic hydrocarbon-degrading *Pseudomonas bauzanensis* Strain W13Z2', *Genome Announcements*, 2: e01049–14
- et al. (2010) 'Cryptic prophages help bacteria cope with adverse environments', *Nature Communications*, 1: 147
- Wang, Z., and Wu, M. (2015) 'An integrated phylogenomic approach toward pinpointing the origin of mitochondria', *Scientific Reports*, 5: 7949
- Westbye, A. B. et al. (2015) 'The gene transfer agent RcGTA contains head spikes needed for binding to the *Rhodobacter capsulatus* polysaccharide cell capsule', *Journal of Molecular Biology*, 428: 477–91
- et al. (2013) 'Phosphate concentration and the putative sensor kinase protein CckA modulate cell lysis and release of the *Rhodobacter capsulatus* gene transfer agent', *Journal of Bacteriology*, 195: 5025–40.
- Williams, K. P., Sobral, B. W., and Dickerman, A. W. (2007) 'A robust species tree for the alphaproteobacteria', *Journal of Bacteriology*, 189: 4578–86
- Yang, Z. (1994) 'Maximum likelihood phylogenetic estimation from DNA sequences with variable rates over sites: approximate methods', *Journal of Molecular Evolution*, 39: 306–14
- Zan, J. et al. (2013) 'The CckA-ChpT-CtrA phosphorelay system is regulated by quorum sensing and controls flagellar motility in the marine sponge symbiont *Ruegeria* sp. KLH11', *PLoS One*, 8: e66346.
- Zhan, Y. et al. (2016) 'A novel roseobacter phage possesses features of podoviruses, siphoviruses, prophages and gene transfer agents', *Scientific Reports*, 6: 30372
- Zhaxybayeva, O., and Doolittle, W. F. (2011) 'Lateral gene transfer', *Current Biology*, 21: R242–6