
Supplementary information

Three families of Asgard archaeal viruses identified in metagenome-assembled genomes

In the format provided by the
authors and unedited

SUPPLEMENTARY INFORMATION

Supplementary text

Enigmatic MGEs of asgardarchaea

Two circular contigs, 7H_11 and 8H_18, were nearly identical and targeted by identical spacers affiliated to Thorarchaeia (Tables S3 and S4), but were assembled from metagenomes originating from samples collected at different depths (59.5 and 68.8 mbsf, respectively). Notably, 7H_42 discovered in our samples from the offshore Shimokita Peninsula, Japan was found to be related to Ga0114923_10000127 and Ga0209976_10000148 originating from the sediment samples from the Sumatra, Indian Ocean (Fig. S1), attesting to the consistency of this emerging group of Asgard MGEs.

The seven MGEs encode diverse proteins involved in DNA replication, repair and metabolism, which are common in MGE and viral genomes but, with the exception of 7H_42, Ga0114923_10000127 and Ga0209976_10000148 which form a group, display little overlap in gene content (see below). Nevertheless, 8H_18, 10H_0 and 7H_42 encode Topo mini-A homologs. Phylogenetic analysis of these proteins showed that, whereas 10H_0 and 7H_42 formed a clade with the Topo mini-A of wyrdvirus WyrdV4, 7H_42 branched among other archaeal sequences (Extended Data 4), suggestive of active exchange of Topo mini-A genes among archaeal viruses and MGEs. The 10H_0 and 7H_42-like MGEs as well as verdandiviruses and wyrdviruses encode multiple non-orthologous Zn-finger proteins, which might be involved in transcription regulation or mediate protein-protein interactions. 10H_0 and 7H_42-like MGEs also share homologs of proliferating cellular nuclear antigen (PCNA) and transcription initiation factor B (TFB) (Fig. S1), both of which have been previously identified in archaeal viruses. For instance, PCNA is encoded by several tailed archaeal viruses infecting halophilic archaea^{1,2} and spindle-shaped viruses of Nitrososphaeria³, whereas TFB homologs are encoded by certain rod-shaped viruses infecting hyperthermophilic Thermoproteota⁴. 10H_0 also encodes Cdc6/Orc1-like origin recognition protein, nucleoside 2-deoxyribosyltransferase, DNA lyase, MazG-like nucleotide pyrophosphohydrolase and bifunctional (p)ppGpp synthase/hydrolase as well as several DNA methyltransferases and nucleases. By contrast, 7H_42-like MGEs encode a DNA primase-superfamily 3 helicase fusion protein that are commonly found in diverse MGEs including diverse varidnaviruses infecting eukaryotes, a Rad51-like recombinase, several nucleases and chromatin-associated proteins containing the HMG domain. The larger elements also encode auxiliary metabolic genes, including PAPS reductase, sulfatase, methylthiotransferase, and enzymes involved in carbohydrate metabolism, which could boost the metabolic activities of the respective hosts. For the smaller contigs, 8H_18 and 8H_67, the vast majority of genes were refractory to functional annotation even using the most sensitive available sequence similarity detection tools, such as HHpred⁵.

Auxiliary gene content of asgardarchaeal viruses

By dsDNA virus standards, genomes of verdandiviruses, skuldviruses and wyrdviruses are relatively small (≤ 20 kb). Thus, the corresponding gene contents are streamlined to include largely the core functions required for virion morphogenesis and genome replication. Nevertheless, some of these viruses encode auxiliary functions, including metabolic genes. In particular, verdandivirus VerdaV1 (and 10H_0 MGE) encode phosphoadenosine phosphosulfate (PAPS) reductase (also known as CysH), an enzyme reducing 3'-phosphoadenylylsulfate to phosphoadenosine-phosphate using thioredoxin as an electron donor. PAPS reductases have been previously identified in certain bacteriophages⁶⁻⁸ and tailed haloarchaeal viruses¹, where they are thought to confer selective advantage to the host cells through facilitating sulfur metabolism and/or synthesis of sulfur-containing amino acids⁷. PAPS reductase of VerdaV1 might perform a similar function.

Wyrdviruses WyrdV2 and WyrdV6 carry a block of three genes coding for dUTPase, thymidylate synthase X (ThyX) and an uncharacterized protein that is conserved in some phages and is annotated as nucleotide modification associated domain 1 protein (PF07659.13, DUF1599) (Fig. 5). This putative operon is likely to be involved in the biosynthesis of thymidylate from dUTP, to increase the pool of nucleotides available for the synthesis of viral DNA. WyrdV3 encodes a homolog of the nucleoside pyrophosphohydrolase MazG, which in bacteria prevents programmed cell death by degrading the central alarmone, ppGpp⁹. MazG is highly conserved in tailed bacteriophages infecting cyanobacteria¹⁰. Biochemical characterization of a cyanophage MazG has shown that, instead of degrading ppGpp, it preferentially hydrolyses dGTP and dCTP¹¹. Thus, MazG homolog in WyrdV3 might either function in disarming antiviral systems triggered by nucleotide-based alarmones, such as ppGpp, or in adjusting the intracellular nucleotide concentrations for optimal viral genome synthesis. Notably,

MazG homologs are also encoded by asgardarchaeal MGEs 10H_0, 7H_42, Ga0114923_10000127 and Ga0209976_10000148 (Fig. S2).

None of the known archaeal viruses encodes its own RNA polymerase¹². Nevertheless, various transcription regulators with HTH, Zn-finger or ribbon-helix-helix domains are abundantly encoded in archaeal virus genomes¹³. This is also the case with asgardarchaeal viruses described herein. Verdandiviruses and wyrdviruses encode multiple non-orthologous Zn-finger proteins, whereas skuldviruses encode several proteins with HTH domains (Fig. 4a). In addition, WyrdV1 (as well as 10H_0, 7H_42, Ga0114923_10000127 and Ga0209976_10000148) encodes a transcription initiation factor B (TFB), a homolog of eukaryotic TFIIB, which guides the initiation of RNA transcription¹⁴. Among archaeal viruses, TFB homologs have been previously identified only in certain rod-shaped viruses infecting hyperthermophilic archaea⁴. Thus, Asgard viruses appear to fully rely on the core transcription machinery of their hosts but encode various transcription factors that could be involved in the recruitment of this machinery for expression of viral genes as well as in the regulation of virus gene transcription. As mentioned above, some of the genes regulated by these transcription factors are likely to encode antidefense proteins.

WyrdV1 and WyrdV3 encode homologs of the carbohydrate-specific 3'-O-methyltransferase¹⁵. In many archaeal viruses, the structural proteins are glycosylated by either the virus or host encoded glycosyltransferases, although the biological role of this post-translational modification remains unclear. The methyltransferase of WyrdV1 and WyrdV3 could participate in modification of the glycans attached to the virion proteins.

Supplementary references

- 1 Mizuno, C. M. *et al.* Novel haloarchaeal viruses from Lake Retba infecting Haloferax and Halorubrum species. *Environ Microbiol* **21**, 2129-2147, doi:10.1111/1462-2920.14604 (2019).
- 2 Raymann, K., Forterre, P., Brochier-Armanet, C. & Gribaldo, S. Global phylogenomic analysis disentangles the complex evolutionary history of DNA replication in archaea. *Genome Biol Evol* **6**, 192-212, doi:10.1093/gbe/evu004 (2014).
- 3 Kim, J. G. *et al.* Spindle-shaped viruses infect marine ammonia-oxidizing thaumarchaea. *Proc Natl Acad Sci U S A* **116**, 15645-15650, doi:10.1073/pnas.1905682116 (2019).
- 4 Baquero, D. P. *et al.* New virus isolates from Italian hydrothermal environments underscore the biogeographic pattern in archaeal virus communities. *ISME J* **14**, 1821-1833, doi:10.1038/s41396-020-0653-z (2020).
- 5 Gabler, F. *et al.* Protein Sequence Analysis Using the MPI Bioinformatics Toolkit. *Curr Protoc Bioinformatics* **72**, e108, doi:10.1002/cpb.108 (2020).
- 6 Mara, P. *et al.* Viral elements and their potential influence on microbial processes along the permanently stratified Cariaco Basin redoxcline. *ISME J* **14**, 3079-3092, doi:10.1038/s41396-020-00739-3 (2020).
- 7 Summer, E. J., Gill, J. J., Upton, C., Gonzalez, C. F. & Young, R. Role of phages in the pathogenesis of Burkholderia, or 'Where are the toxin genes in Burkholderia phages?'. *Curr Opin Microbiol* **10**, 410-417, doi:10.1016/j.mib.2007.05.016 (2007).
- 8 Farlow, J. *et al.* Genomic characterization of three novel Basilisk-like phages infecting *Bacillus anthracis*. *BMC Genomics* **19**, 685, doi:10.1186/s12864-018-5056-4 (2018).
- 9 Gross, M., Marianovsky, I. & Glaser, G. MazG -- a regulator of programmed cell death in *Escherichia coli*. *Mol Microbiol* **59**, 590-601, doi:10.1111/j.1365-2958.2005.04956.x (2006).
- 10 Sullivan, M. B. *et al.* Genomic analysis of oceanic cyanobacterial myoviruses compared with T4-like myoviruses from diverse hosts and environments. *Environ Microbiol* **12**, 3035-3056, doi:10.1111/j.1462-2920.2010.02280.x (2010).
- 11 Rihtman, B. *et al.* Cyanophage MazG is a pyrophosphohydrolase but unable to hydrolyse magic spot nucleotides. *Environ Microbiol Rep* **11**, 448-455, doi:10.1111/1758-2229.12741 (2019).
- 12 Krupovic, M., Cvirkait-Krupovic, V., Iranzo, J., Prangishvili, D. & Koonin, E. V. Viruses of archaea: Structural, functional, environmental and evolutionary genomics. *Virus Res* **244**, 181-193, doi:10.1016/j.virusres.2017.11.025 (2018).
- 13 Iranzo, J., Koonin, E. V., Prangishvili, D. & Krupovic, M. Bipartite network analysis of the archaeal virosphere: Evolutionary connections between viruses and capsidless mobile elements. *J Virol* **90**, 11043-11055, doi:10.1128/JVI.01622-16 (2016).
- 14 Werner, F. & Grohmann, D. Evolution of multisubunit RNA polymerases in the three domains of life. *Nat Rev Microbiol* **9**, 85-98, doi:10.1038/nrmicro2507 (2011).
- 15 Bernard, S. M. *et al.* Structural basis of substrate specificity and regiochemistry in the MycF/TyfF family of sugar O-methyltransferases. *ACS Chem Biol* **10**, 1340-1351, doi:10.1021/cb5009348 (2015).

Supplementary figures

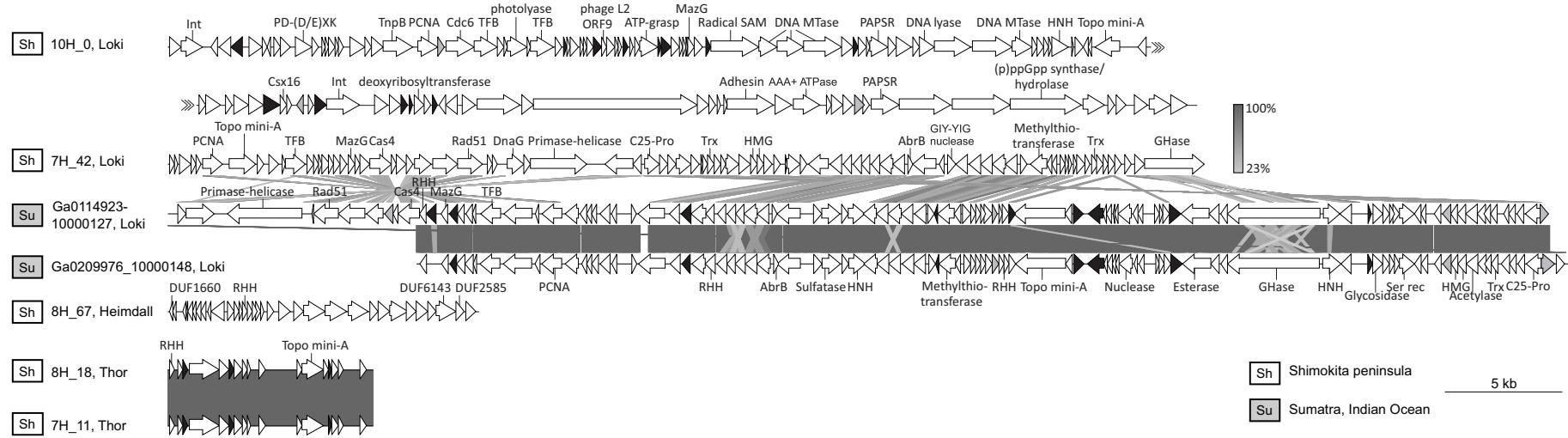


Figure S1. Genome maps of asgardarchaeal CRISPR-targeted MGEs. Due to the lack of known virus hallmark genes, the MGEs are not identified as viruses, but could represent novel virus groups. Genes encoding putative DNA-binding proteins with Zn-binding and helix-turn-helix domains are colored in black and grey, respectively. Grey shading connects genes displaying sequence similarity at the protein level, with the percent of sequence identity depicted with different shades of grey. Abbreviations: Int, integrase; PD-(D/E)XK, PD-(D/E)XK family nuclease; PCNA, proliferating cellular nuclear antigen; TFB, transcription initiation factor B; MTase, methyltransferase; HNH, HNH family nuclease; PAPSR, phosphoadenosine phosphosulfate reductase; C25-Pro, C25-family protease; Trx, thioredoxin; GHase, glycoside hydrolase; RHH, ribbon-helix-helix domain-containing DNA binding protein; Ser rec, serine superfamily recombinase; HMG, high mobility group domain-containing chromatin-associated protein.

a

>P17312 TERL_BPT4 Terminase, large subunit OS=Enterobacteria phage T4 OX=10665 GN=17 PE=1 SV=1
Probab=100.00 E-value=4.e-34 Score=301.28 Aligned_cols=403 Identities=15% Similarity=0.129 Sum_probs=304.0 Template_Neff=10.300

b

>A0A0U5AF03 CAPSD_BPK22 Major capsid protein OS=Pseudomonas phage KPP22 OX=1772250 PE=1 SV=1
Probab=99.97 E-value=3.5e-29 Score=244.70 Aligned_cols=275 Identities=1% Similarity=0.022 Sum_probs=0.0 Template_Neff=7.800

Q ss_pred	CccccCCCCHHHHHHHHHHHhhccchhccccCCCCCCCCC-ccceEEeeeeeeccecc--ccCccccccccc-cceEEEE
Q VerdaV1	2 VEIHGVSERQWKIDTQAQLMRPEIVYTRPLPSMQIG-KGKNAFKYFQATDIFTYVN--EGNRFSYEDKRT-GARKIRDA
Q Consensus	2 ~~~~~~l~~~~~w~~~Id~~v~~~~~L~ar~p~~~~~G~~~~~d~~~~~g~~~~~v~~~~~v~~~~~v~~~~~v
T Consensus	73 ~~~~~ip~~~~~id~~vi~~~~~l~~~~~li~~~~~g~~~~~w~~~~~d~~~~~G~~~~~d~~~~~p~~~~~v~~~~~v
T AAOU5UF03	73 TPSIPTIQFLQTLPGLVKVMTAARKIDTEIG-IDTVGSWEQEIVQGIVPAGTAVEGYGDHTN1PLTSWNANFERTTI
T ss_pred	CCCCchHHHHHHHhhccchhcccccccCHHHcC-ccccCCcceeEEEEEEecCceEEecCCCCcccceeEeeeeeee

Q ss_pred	CCCCCCCCCChhhceecchh--cCCCCCCCeeEEEEECC-----CeEEEEecCCeEcCCCCCeeee	
Q VerDaV1	224 GEKGTEPFIPDINQITEDAL---YNGTLSTATQAFLLWKWDI-----NCNYIAEAYPIHRAGIPKNKEFEGDID	289 (318)
Q Consensus	224 ~~~~~~i~~~~~i~~~~~d~i~~~~~i~~~~~i~~~~~i~~~~~i~~~~~i~~~~~i~~~~~i~~~~~i	289 (318)
T Consensus	.+ . ..+++++. +.+ .+ .+++++.++ +++.+.+ +++.+++.++. ++	
T AOA0U5AF03	292 ---p~i~----pel~a~~g~-----g~-----l~~~~~g~~~~~v~~~~~v~~~~~v~~~~~v	360 (382)
T ss_pred	292 ----KMRV~IISPAPELSGVQMKQA---PEDALVFLFVEDVNAAVDGSTDGGSVFSVLOSVKFITLGVEKRAKSYVED	360 (382)

Q ss_pred	EEEEEecEEecccCCcEEE	
Q VerlaV1	290 YALRWAGCFLPKNPKGAVYV	309 (318)
Q Consensus	290 ~~~r~gGv~~~~~	309 (318)
	..++ ++ ++ .	
T Consensus	361 ~~~tG~~~i~~rP~ai~~~	379 (382)
T AOA005AF03	361 FSNGTAGALC-KRPWAVVRY	379 (382)
T ss_pred	eeccccEEE-ccccchhh	

Figure S2. Results of the HHsearch analysis queried with the putative large terminase subunit (a) and the major capsid proteins (b) of verdandivirus VerdaV1. H(h), α -helix; E(e), β -strand; C(c), coil.

a

>P15794 CAPSD_BPPM2 Major capsid protein P2 OS=Pseudoalteromonas phage PM2 OX=10661 GN=II PE=1 SV=2
Probab=98.32 E-value=0.0012 Score=61.40 Aligned_cols=219 Identities=12% Similarity=0.056 Sum_probs=0.0 Template_Neff=8.200

Q ss_pred	eeecCeeeeeeccccEEEEEccCcccCchhhheeecccccEEEEECCccccHHHHHHHHHHHH	
Q SkulDv1	15 AESENKHVKIDNTFPPIKKVIIHNNPGDFAGGSTGVLIANTIDTLNRLRVFHGKEVISL-IGDVDTAPFAQLMREANKL	93 (277)
Q Consensus	15 ~~~~~~ti~tf~ix~I~~~~ga~sGgatG~Va~N~a~s~rv~nGk~I~~-DG~~~vd~s~gi~lRRE~~~~~ +.+.++ + -.-..+ .+..+ .+..-+ . + .+..+ + ..+ ++	93 (277)
T Consensus	14 ~~g~at~aa~lp~g~t~Y~~i~l~t~-----~a~i~I~v~NG~i~i~v~-----~L~n~n~y	71 (269)
T P15794	14 AAGNSCSIKLPIGTYEVIDLRYS-----GVTPSQKINVRVELDGRLLSTYKTL-----NDLILENTR	71 (269)
T ss_pred	CCCCEEEECCCCeeEEEEEEEc-----CCHHHHCeEEEEBCCCCCeCCH-----HHHHHHHHHH	

Q ss_pred	hcccC CCCCCcEEEcCCcCCCC-----CeEEEEEEecceeecccCCCCceeeeeeccccCCCCe		
Q Skuldv1	94 R NKVADADEYFIDPAKIA GRG-----HNAYIDFRNFTIANMNDGR TTYGATIDMVEVGTPGR	154	(277)
Q Consensus	94 ~~VA~~~d~y~I~FP~aip~g-----dvgi~l~~~tsqa~gGdr~T~at~D~i~v~~~Kg~ + + + . + + + . . + - + + . + . + + + , + + + . + + + + . + ..	154	(277)
T Consensus	72 ~g~~~~~g~l~i~F~~~~~a~~~~~al~t~~~s~l~ve~i~a~~~~~p~L~w~a~~~~~	140	(269)
T P15794	72 H- KRKIKAGVVSFHVFVRPEMKGVNTDLVQQRMFALGTVGTLTCE1KF DIDEAA-----GPKL SIAQKSVGT A	140	(269)
T ss_pred	h-c CCCCCcEEEcCChhcCCCChHheeehhccCCCCcEEE EEEccccCCCC-----CCEEEE EeCCCCC		

Q ss_pred	ceeeccceEEEEEcC--cC-----CHHHeEEEEEcCcCcCeEEEEEecccC		
Q Skuldv1	232 FS1ALGTGFVMP1PKQ--AI-----SASQLKLISQIDSAGTNYEIVWMLTSK	277	(277)
Q Consensus	232 s~vA-s-G~~~i~fp---ki-----s~tlkl~~~~tAg~e~h~~~~~	277	(277)
	.++.-+ +++.+=.=... .++ + ++..++ ++=--.-.+..++		
T Consensus	211 ~~~~pqag~~~~DF~~~~g~~~~~a~~~~~a~~~~~a~~~~~i~~~~~E~~~~~	265	(269)
T P15794	211 GKAV-LDNTYTDIFMLEGDYVQSLLDQM1QDLRKLDSTMDEQAKIIVEYGMWS	265	(269)
T ss_pred	Cccc-CCeEEEEECCCCChhhhhhccccceEEEEEeCcCcCcEEEEEecccc		

b

>3J31_B Major capsid protein; Double jelly roll fold; 4.5A; Sulfolobus turreted icosahedral virus
Probab=97.84 E-value=0.013 Score=58.95 Aligned_cols=242 Identities=13% Similarity=0.060 Sum_probs=0.0 Template_Neff=6.300

Figure 1. The effect of the number of hidden neurons on the performance of the proposed model.

Q	ss	pred	-EEEEECCccccCCCCHHHHHHHHHCCc-eccccEEEEEccC-----c--CHHHeEEEEccCccceEEE		
Q	SkuldfV1		207 -LELKGGSNIIRREGSPLQLEKTGKFS-IALGFTFWIPIKQ-----A- ISASKLQKLSIQSDAGTNIEVH	270	(277)
Q	Consensus		207 -L~Is~g~~~i~DGsikl~~~~~Ks~`~Va~s~G~~~i~fp-----k-is~tLkL~s~t~tag~e~h	270	(277)
T		+...+.++...- +++.+.+-..+ +++- -+ ..-++!+.++!.---		
T	Consensus		237 ~~~~~~l~~~~~l~~~~~l~~~~~g~~~~~idF~~~~~g~l~~~~~L~~~~~g~~~~~vi	312	(345)
T	3J31_B		237 KIVRGVPT-DKIKVSWAALQAENQAEYQVAPYSGASAIDFRKYFNGDLDTAHPDSIEYDIALQNQDNVLYS	312	(345)
T	ss_dssp		BECSSCE-----EECBBBHHHHHHCCSCSSTTTEEEGGTTTTEEESSCCSSCCSEEEEEE		
T	ss	pred	EeeCCC-----eeeeCHHHHHHHHHCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC		

C

>P22535 CAPSD_BPPRD Major capsid protein P3 OS=Enterobacteria phage PRD1 OX=10658 GN=III PE=1 SV=2
Probab=96.16 E-value=1.5 Score=46.49 Aligned_cols=247 Identities=11% Similarity=0.091 Sum_probs=0.0 Template_Neff=5.60

Q ss_pred	HHHHHHHHHHccCcCc	ccccEEccCCCC	-eEEEE	
Q SkuldiV1	84 IQMREANKLNRKVADA	-DEYFEIDFPKAIRGHH-	-NAYID	121 (277)
Q Consensus	84 i~I~RE~P~N~A~V~A~	d~y~`I~FP~a~ip~g~	d~gvq~`	121 (277)
	+.+++++....	.~.~+ ++ ++ ++.	++++	
T Consensus	117 --L~--n~-----	~-as~Ag~t~v~f~l~i~Pval~--D~G~il~qn~t~l~L~	193 (395)	
T P22535	117 --LHFVNNTAKQGAPFLSSMVTDSPPIKYGVDMNNVIDAPATIAGATGELTMYWVPLAYSETDLTGAVLANPVQSQRKLK	193 (395)		
T ss_pred	---HHHHHHHccCccc==ccccCCCCccccCCCCccccCccCeEEEEE	==ccccCCCCccccccccCCCCccccCccCeEEEEE		

Q ss_pred	---	---CCHHHEEEEEEccCcceEEEEee	
Q SkuldV1	248	--AISASQSLKLISQDAGTNIEVHMLT	274 (277)
Q Consensus	248	--kis~+tLk~~~~+tAgt~~e~h~~~	274 (277)
		----+. .+.-+. .+.-+ .	
T Consensus	349	gi-t~~~~~l~~~~~	378 (395)
T P22535	349	PIYTLLQYGNVGFVNPKTVNQNARLLMGYE	378 (395)
T ss_pred		cccccccccEEEEEccccccccceEEEEee	

Figure S3. HHsearch profile-profile comparisons between the putative major capsid protein of skuldvirus SkuldV1 and the double jelly-roll major capsid proteins of (a) corticovirus PM2, (b) turrivirus STIV, and (c) tectivirus PRD1. H(h), α -helix; E(e), β -strand; C(c), coil.

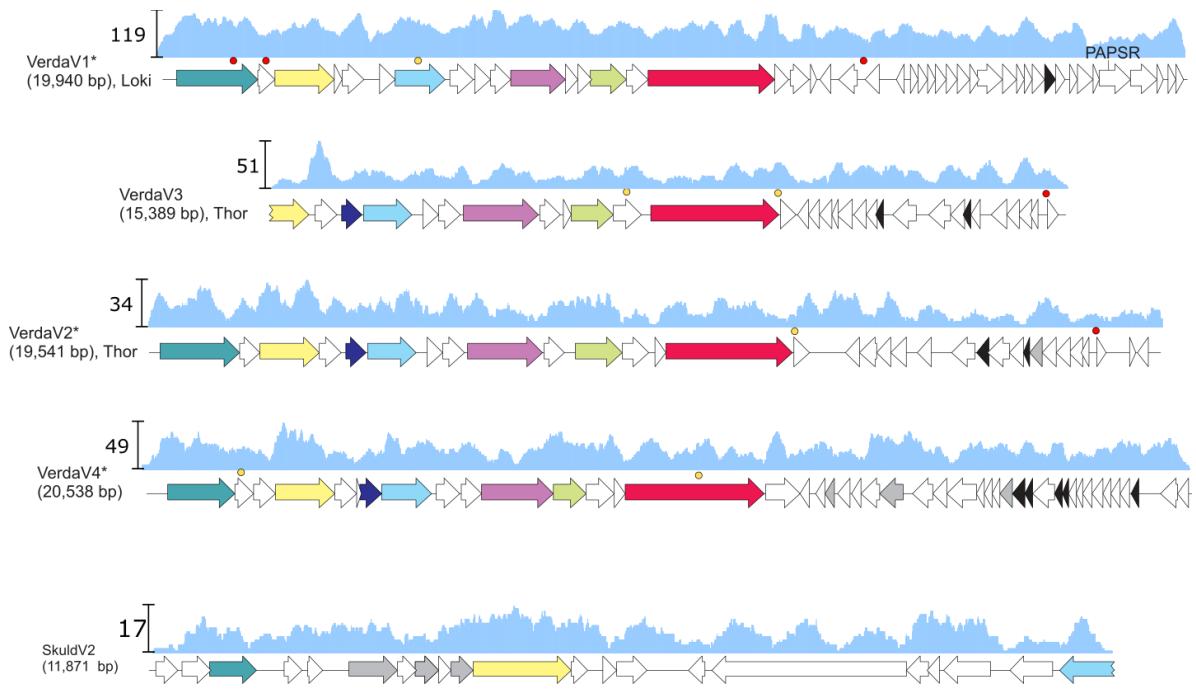


Figure S4. Read depth along viral contigs. The scale on the left shows the maximal coverage for each contig. The open reading frames are colored the same way as in the corresponding Figures 3 and 4 for verdandiviruses and skuldviruses, respectively. Colored circles represent the positions of protospacers targeted by asgardarchaeal CRISPR arrays.

Supplementary Data 1. Output of CRISPRDetect output for CRISPR arrays analyzed in this work.