

Logan: Planetary-Scale Genome Assembly Surveys Life's Diversity

Rayan Chikhi¹, Téo Lemane², Raphaël Loll-Krippleber^{3,4}, Mercè Montoliu-Nerin⁵, Brice Raffestin¹, Antonio Pedro Camargo^{6,7}, Carson J. Miller⁸, Mateus Bernabe Fiampenghi⁷, Daniel Paiva Agustinho⁹, Sina Majidian¹⁰, Greg Autric¹¹, Maxime Hugues¹², Junkyoung Lee¹³, Roland Faure¹, Kristen D. Curry¹, Jorge A. Moura de Sousa¹, Eduardo P. C. Rocha¹, David Koslicki¹⁴, Paul Medvedev¹⁴, Purav Gupta^{4,15}, Jessica Shen^{4,15}, Alejandro Morales-Tapia^{4,15}, Kate Sihuta^{3,4}, Peter J. Roy^{3,4}, Grant W. Brown^{3,4}, Robert C. Edgar¹⁶, Anton Korobeynikov¹⁶, Martin Steinegger¹², Caleb A. Lareau¹⁷, Pierre Peterlongo¹⁸, and Artem Babaian^{4,15}

¹Institut Pasteur, Université Paris Cité, CNRS UMR3525, Paris, France

²Génomique Métabolique, Genoscope, Institut de Biologie François Jacob, CEA, CNRS, Univ. Evry, Université Paris-Saclay, Evry, France

³Department of Biochemistry, University of Toronto, Toronto, Canada

⁴The Donnelly Centre for Cellular & Biomolecular Research, University of Toronto, Toronto, Canada

⁵Tree of Life, Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton CB10 1SA, UK

⁶Department of Biochemistry, Institute of Chemistry, University of São Paulo, São Paulo, SP, Brazil

⁷DOE Joint Genome Institute, Lawrence Berkeley National Laboratory, Berkeley, CA, USA

⁸Department of Microbiology, University of Washington, Seattle, WA, USA

⁹Human Genome Sequencing Center, Baylor College of Medicine, Houston, Texas, USA

¹⁰Department of Computer Science, Johns Hopkins University, 3400 North Charles St., Baltimore, MD 21218, USA

¹¹Amazon Web Services, Paris, France

¹²Amazon Web Services Inc., Seattle, USA

¹³School of Biological Sciences, Seoul National University, South Korea

¹⁴Computer Science and Engineering, Biology, and the Huck Institute of the Life Sciences, Pennsylvania State University, University Park, PA, USA

¹⁵Department of Molecular Genetics, University of Toronto, Toronto, Canada

¹⁶Independent

¹⁷Computational and Systems Biology Program, Memorial Sloan Kettering Cancer Center, New York, NY 10065, USA

¹⁸University Rennes, Inria, CNRS, IRISA - UMR 6074, 35000 Rennes, France

Abstract

The breadth of life's diversity is unfathomable, but public nucleic acid sequencing data offers a window into the dispersion and evolution of genetic diversity across Earth. However the rapid growth and accumulation of sequence data have outpaced efficient analysis capabilities. The largest collection of freely available sequencing data is the Sequence Read Archive (SRA), comprising 27.3 million datasets or 5×10^{16} basepairs. To realize the potential of the SRA, we constructed Logan, a massive sequence assembly transforming short reads into long contigs and compressing the data over 100-fold, enabling highly efficient petabase-scale analysis. We created Logan-Search, a k -mer index of Logan for free planetary-scale sequence search, returning matches in minutes. We used Logan contigs to identify >200 million plastic-degrading enzyme homologs, and validate novel enzymes with catalytic activities exceeding current reference standards. Further, we vastly expand the known diversity of proteins (30-fold over UniRef50), plasmids (22-fold over PLSDB), P4 satellites (4.5-fold), and the recently described Obelisk RNA elements (3.7-fold). Logan also enables ecological and biomedical data mining, such as global tracking of antimicrobial resistance genes and the characterization of viral reactivation across millions of human BioSamples. By transforming the SRA, Logan democratizes access to the world's public genetic data and opens frontiers in biotechnology, molecular ecology, and global health.

47 1 Main

48 DNA sequencing has revolutionized our perspective on life's diversity, yet the majority of the world's
49 sequencing data are inaccessible to systematic search and analysis. The Sequence Read Archive (SRA)
50 houses over 50 petabases (Pbp; 5.0×10^{16}) of public sequencing data, and is growing exponentially (Fig. 1a)
51 [1]. This data represents billions of dollars of global research output, spanning all known life and covering
52 every continent (Fig. 1).

53 Analyses of the SRA have yielded profound scientific discoveries, from hundreds of thousands of novel
54 viruses to shifts in antibiotic resistance patterns [2, 3, 4, 5, 6, 7, 8, 9]. Yet the methods for massive-scale
55 sequence analyses, based on assembly or k -mer indexing, face computational and economic constraints.
56 The largest collection of public assembled data, NCBI GenBank WGS, spans less than 3% of the SRA [10],
57 while k -mer indexing of reads have not scaled beyond 9% of direct SRA data (Table 1).

58 To enable SRA-wide analysis, we developed Logan, an assembly of 96% of the SRA (27 million acces-
59 sions as of December 2023) and a suite of associated tools. Using massively parallel cloud processing, we
60 transformed 44.1 petabases of raw sequencing data into 0.9 petabases of long assembled contigs. Logan
61 achieved a >100-fold compression of the original SRA data, and is 36-fold larger than GenBank WGS.
62 Logan assembly fundamentally changes the economics and speed of SRA-wide searches. To demonstrate
63 this, we developed Logan-Search, a k -mer index for finding a nucleotide sequence across all Logan as-
64 semblies in minutes. For protein homology search, we aligned using DIAMOND2 all Logan contigs to a
65 protein query in 11 wall-clock hours, with a near 20-fold cost decrease relative to previous methods such
66 as Serratus [11, 2].

67 We demonstrate Logan's utility by making discoveries in three domains. First, for bioprospecting,
68 we performed an SRA-wide search for homologs of plastic-degrading enzymes with sensitivity to $\approx 40\%$
69 amino acid sequence identity, discovering over 200M novel enzymes, including several which we validated
70 as having higher and more varied catalytic activities than previous standards. Second, for scalable clinical
71 discovery, we used Logan-Search to screen millions of human datasets for viral gene expression sequences.
72 We uncovered recurrent Human Herpesvirus-6 reactivation in tumor-infiltrating lymphocyte therapy prod-
73 ucts, broadening the characterization of viral reactivation in cell therapies *ex vivo* [12]. Third, we mined
74 Logan contigs for proteins, plasmids and subviral elements. This effort massively expanded the known
75 protein universe, with a 30-fold increase in protein diversity over UniRef at 50% amino-acid clustering
76 identity. We also obtained a 22-fold increase in plasmid families, a 4.5-fold increase in the diversity of P4
77 satellites, and a 3.7-fold expansion of the recently described Obelisk-like species [9].

78 A Petabyte-Scale Assembly and Search Engine of the Sequence Read Archive

79 To construct the Logan assemblage we applied two complementary assembly strategies to the entirety of
80 the SRA (as of 2023-12-10), comprising 27.3 million SRA accessions (Fig. 1b). The first strategy generates
81 unitigs, which are near-lossless representations aiming to preserve sequence content from a sample, making
82 them ideal for sensitive k -mer-based search. The second strategy builds on the unitigs to create contigs,
83 which form longer, consensus sequences by resolving small biological variations. Contigs are optimized
84 for protein identification and other downstream analyses.

85 In total, the Logan assemblage is 4.59 Ppb of unitigs (2.18 petabytes compressed) and 0.90 Ppb of
86 contigs (0.31 petabytes compressed, Fig. 1b), generated in approximately 30 hours of wall-clock time, using
87 a peak of 2.18 million CPU cores (Methods, Extended Data Fig 1). All assemblies and documentation are
88 publicly hosted and freely available (<s3://logan-pub/>, <https://github.com/IndexThePlanet/Logan>).

89 To make the Logan assemblage rapidly searchable we developed Logan-Search, a 1 petabyte k -mer
90 ($k=31$) index of Logan unitigs from 23.4M SRA accessions. For queries of up to 1 kb, Logan-Search returns
91 all indexed SRA accessions containing a user-set fraction of the query k -mers, along with associated SRA
92 metadata. This provides rapid insights into the environmental and geographic distributions of a gene
93 (Fig. 1c). Logan-Search surpasses previous non-Logan based SRA sequence-search efforts [8, 4, 6] by at

94 least an order of magnitude in both the number of accessions and total coverage (Table 1). To democratize
95 SRA-scale sequence exploration, we deployed Logan-Search as a web interface (<https://logan-search.org>), enabling researchers to freely query Logan unitigs in a few minutes.
96

97 Logan's 0.9 petabases of contigs span the entire tree of life, providing orders of magnitude more
98 assembled data for nearly every sequenced species (Fig. 1d). For key model organisms and agricultural
99 species, the amount of assembled data increases dramatically, rising from less than 4 Tbp in GenBank
100 to approximately 285 Tbp for *Homo sapiens* (70-fold increase), from 0.9 Tbp to 21 Tbp for *Bos taurus*
101 (cattle, 23-fold increase), from 0.5 Tbp to 8.8 Tbp for *Gallus gallus* (chicken, 18-fold increase), and from
102 0.2 Tbp to 7 Tbp for *Zea mays* (maize, 35-fold increase).

103 The assembled metagenome samples in Logan exceed 130 terabases (over 5.7 million accessions), 144-
104 fold more bases than a comparable collection of short read assembled metagenomes (MGNify [13]). To
105 quantify the novel sequence content within Logan's metagenomes, we used sketching to estimate the
106 number of distinct 31-mers and compared this to the entirety of the GenBank WGS database. Logan's
107 metagenome collection contains an estimated 33.9 trillion distinct *k*-mers that appear two or more times
108 in an accession, a nearly 4-fold increase over the 8.7 trillion found in all of GenBank WGS. Crucially, of
109 Logan's 33.9 trillion metagenomic *k*-mers, 32.3 trillion (95%) are not present anywhere in the GenBank
110 WGS database, highlighting a vast reservoir of previously uncharacterized genetic diversity.

111 Expanding the arsenal of plastic-active enzymes

112 Logan is a powerful tool for enzyme discovery. Since the mid-20th century humans have manufactured in
113 excess of 12 gigatons of plastics, with the majority ending up as waste which degrades into micro- and
114 nanoplastics. These particles infiltrate global ecosystems and food supplies, where they bio-accumulate
115 to high incidence in humans [14, 15]. The 2016 discovery of a polyethylene terephthalate (PET) plastic
116 degrading enzyme in *Ideonella sakaiensis* (*IsPETase*) has catalyzed bio-prospecting and bioengineering
117 efforts to identify novel and high-efficiency enzymes for recycling and remediating plastic polymers [16,
118 17, 18, 19, 20].

119 To expand the diversity of plastic-active enzymes, we created a search query from the 213 validated
120 plastic-active enzyme sequences in the PAZy database [18]. These sequences group into 11 CATH protein
121 domains [21], showing polyphyletic activity against varying plastic substrates (Fig. 2a). First, we searched
122 the 213 sequences in the NCBI nr database (0.003 Pbp, DIAMOND2 blastp) and recovered 2.73 million
123 distinct sequences or 1.05 million non-redundant enzyme homologs (clustered at 90% amino acid identity,
124 aaid).

125 To expand this, we then queried these 1.05 million enzymes across Logan contigs, recovering 1.12
126 billion matching sequences. From these, 385 million sequences (34.3%) matched an nr enzyme at 90%
127 identity. We clustered the remaining 735 million sequences (65.7%) into 215.7 million non-redundant
128 enzyme homologs (Fig. 2b). Overall, Logan provided a 205-fold expansion of sequence diversity over NCBI
129 nr, spanning multiple plastic-active domains, including a 190-fold expansion of A/B hydrolases, of which
130 *IsPETase* is a member (Fig. 2c). This dataset, *PETadex*, represents the most comprehensive and diverse
131 collection of candidate plastic active enzymes. To facilitate the development of global research solutions
132 for plastic remediation, we are releasing the *PETadex* data freely, and without restriction (<https://github.com/ababaian/petadex>).

133 To test if *PETadex* candidate plastic-active enzymes contain catalytically active sequences, we developed
134 a quantitative high-throughput enzyme screen using a PET subunit substrate, bis(2-hydroxyethyl)
135 terephthalate (BHET). Candidate PETase enzymes were expressed as cell surface displayed or secreted
136 constructs in baker's yeast, *Saccharomyces cerevisiae*. PET conversion activity was measured via a colo-
137 metric halo around yeast colonies, which corresponds to the formation of 2-hydroxyethyl terephthalate
138 (MHET) or a higher molecular weight "halo product" (Fig. 2d, Extended Data Fig. 3).

139 As an initial screen, we selected full-length *IsPETase* or PAZy A/B hydrolase-like homologs (40%
140 aaid). From 2,272 unique matches, we synthesized 161 randomly selected enzymes, and 21 ancestral

142 reconstructions (AR). We screened these enzymes over six timepoints and in the two expression systems,
143 which revealed 35/161 (22%) of the natural, and 8/21 (38%) of the AR sequences had plastic-activity
144 (Fig. 2f, Extended Data Fig. 4a). The most active enzymes identified showed overall conversion compara-
145 ble to *IsPETase* with varying preferences for MHET or halo product formation (Extended Data Fig. 4a).
146 Inspection of the enzyme phylogeny revealed a clade enriched for halo-formation activity, which was
147 re-sampled for an additional 13 enzymes. Resampling yielded 9/13 (69%) plastic-active enzymes. High-
148 performance liquid chromatography (HPLC) assessment of BHET conversion activity of the two most
149 active enzymes (SRR23008605_28430 and SRR10663367_452477) revealed that these enzymes exceeded
150 *IsPETase* activity for MHET or halo product formation (Extended Data Fig. 4b). Interestingly, these
151 two enzymes also produced the PET monomer product TPA at substantially higher rates than *IsPETase*
152 and 4-fold more TPA than the engineered FAST-PETase (Fig. 2g).

153 Microplastics are projected to increase exponentially and biocatalysts are a means by which these
154 environmental contaminants can be remediated. To address this we created *PETadex*, a free and un-
155 restricted resource of candidate plastic-active enzymes, two orders of magnitude more expansive than
156 previously available. Logan resources such as this enable the deep exploration of the evolutionary land-
157 scape of proteins including identifying candidate enzymes with higher application-specific activities, more
158 complete product yields (TPA formation), or novel chemical functions (halo product formation).

159 Characterization of Human Herpes Virus 6 reactivation in heterogeneous *ex vivo* 160 cultures

161 Our recent work has demonstrated that retrospective assembly and quantification of viral nucleic acids at
162 the petabase scale could reveal new associations between humans and viruses, including in clinical contexts
163 [2],[22]. Specifically, comprehensive mining of Serratus [2] led to the discovery of Human Herpesvirus 6
164 (HHV-6) reactivation in chimeric antigen receptor (CAR) T cells [22], a finding that contributed to revised
165 Food and Drug Administration guidelines requiring the screening of viral reactivation in allogeneic CAR
166 T cells [23].

167 We hypothesized that Logan could further identify instances of viral reactivation in human cells and
168 tissues where viral expression was not considered in the primary analyses. Among the 103 HHV-6 type B
169 (HHV-6B) genes, we selected two transcripts U83 and U91 for sequence query based on high expression,
170 short length (less than 1kb), and known gene function (Extended Data Fig. 6a; Methods). We queried
171 these two transcripts using Logan-Search against 1,476,236 human RNA-sequencing datasets (Fig. 3a).
172 Each query took less than five minutes to complete, resulting in 13 distinct BioProjects with >50% *k*-mer
173 coverage across both HHV-6 transcripts (Fig. 3b; Methods), four of which had known HHV-6 expression.
174 These four projects served as positive controls, which included CD4+ memory T cell cultures annotated
175 by Serratus [24],[25] and CAR T products with previously characterized HHV-6 reactivation [22], [26]).

176 Supported by the recovery of these positive controls, we then considered the nine novel BioProjects.
177 According to the BioSample meta-data, these samples comprised additional gastrointestinal tumors [27]
178 and CAR T cells [22], consistent with previous characterizations in other settings (Fig. 3c). The novel
179 CAR T BioSamples were from a study profiling infusion products and longitudinal profiles from 26
180 patients with B cell acute lymphoblastic leukemia receiving anti-CD19/CD22 CARs [28]. While this
181 study focused predominantly on CAR-intrinsic gene expression changes associated with variable patient
182 outcomes, Logan enables retrospective discovery of HHV-6 reactivation in these samples. Further, as
183 this study was published after the completion of Serratus [2], HHV-6 reactivation in this cohort was
184 not previously annotated or reported by the original authors. These results further supports our prior
185 conclusions of HHV-6 reactivation occurring agnostic of disease or CAR target.

186 Next, we focused on the tumor-infiltrating lymphocytes (TILs) or organoid model annotated BioSam-
187 ples. To the best of our knowledge, no prior reports of viral reactivation had been previously reported
188 in these settings. However, since these systems involve extended cultures of heterogeneous mixtures that
189 include CD4+ T cells, we reasoned that the Logan associations could reflect HHV-6 reactivation in *ex vivo*

190 cell culture settings of T cells [22]. Indeed, in the lung organoid sample [29], we identified a population of
191 HHV-6 super-expressor cells among the CD4+ proliferating T cells (Extended Data Fig. 6b), supporting
192 our prior characterization of a rare cell state responsible for seeding lytic virus following reactivation
193 [22]. The annotated HHV-6+ TIL samples were infusion products profiled from patients with metastatic
194 melanoma treated from three clinical trials [30]. Analysis of the 16 donors profiled with RNA-seq demon-
195 strated high-confidence HHV-6 detection in 5 infusion products, predominantly from the CD4+ sorted
196 populations (Fig. 3d). At least one positive donor was observed in each trial, underscoring that HHV-6
197 reactivation in these adoptive cell therapies is a recurrent phenomenon.

198 As viral RNA accumulation coincides with increased viral DNA copy number [22], we examined
199 additional profiles of clinical TIL products analyzed via chromatin immunoprecipitation and sequencing
200 (ChIP-seq) for a pan-H3 acetylation modification that marks transcriptionally active chromatin. Among
201 19 donors spanning the same three clinical trials, we observed viral reactivation in samples from all three
202 trials, including high HHV-6 expression from a donor (D10) where RNA-seq was not obtained (Fig. 3e;
203 Extended Data Fig. 6c). Further analyses of viral single-nucleotide polymorphisms revealed 72 mutations
204 specific to either donor, excluding the possibility of a common source of HHV-6 contamination during
205 library preparation (Extended Data Fig. 6d). Across both modalities, our results suggest that HHV-6
206 reactivation in T cell therapies occurs independent of exogenous DNA and further implicates the rapid
207 proliferation of T cells *ex vivo* as a critical signal underlying HHV-6 reactivation *in vitro*.

208 Taken together, our analyses shows that Logan effectively uncovers novel biological associations of
209 viruses using existing human genomic profiles. In particular, this vignette reveals that latent HHV-6
210 can reactivate in rare proliferating CD4+ T cells from heterogeneous cell culture conditions, spanning
211 from organoids to adoptive cell therapies with or without genetic engineering. As culture duration is
212 a key determinant of viral reactivation in CAR T cells [22], our observation of HHV-6 reactivation in
213 TIL therapies is consistent with the longer culture durations than widely-used autologous CAR T cell
214 therapies. Our characterization of viral reactivation *ex vivo* motivates further work into gene editing
215 and/or small molecule approaches that can mitigate reactivation to maximize the safety and efficacy of
216 cell therapies [23]. More generally, our results motivate future work to monitor viral reactivation across
217 current and future cell therapies using comprehensive genomics profiling and scalable analyses enabled
218 by Logan.

219 Expanding the Known Universe of Proteins, Plasmids, and Viral Elements

220 Next, we mined the 0.9 petabases of Logan's assembled contigs to reveal order-of-magnitude expansions
221 in the known diversity of proteins and mobile genetic elements. These planetary-scale deep homology
222 searches were completed in as little as 11 hours, using cloud-deployed translated protein-to-nucleotide
223 alignment.

224 **Billions of diverse Logan proteins** Logan expands the known protein universe, with 109.4 billion
225 proteins clustered into 3.0 billion non-redundant sequences at 90% amino acid identity and 90% alignment
226 overlap (Fig. 4d). This represents over an order of magnitude greater set of protein diversity relative
227 to large-scale commercial (BaseData [31]) or public (OMG [32], MGnify [13], BFD [33]) metagenomic
228 resources, and a nearly 30-fold increase over UniRef50.

229 This expanded diversity provides an invaluable resource for downstream applications. In a case study
230 of 100 viral proteins, sensitive searches against a clustering of Logan proteins at 50% amino acid iden-
231 tity produced substantially more diverse multiple sequence alignments compared to searching against a
232 standard database [33]. We observed an approximately 2-fold increase in the Number of Effective Se-
233 quences (Neff: 2.19 to 4.89; Extended Data Fig.5B, left panel). Doubling of Neff values translates into
234 improved protein structure predictions, with mean predicted Local Distance Difference Test (pLDDT)
235 scores rising from 46.7 ("very low") to 88.6 ("high"), and 90 out of 100 proteins modeled at high qual-
236 ity (Extended Data Fig.5B, right panel), highlighting the impact of our expanded protein clustering on

237 improving AI-structure prediction accuracy.

238 **Obelisks** Obelisks are a newly discovered clade of viroid-like agents. Initially two complete, circular
239 species were identified as persistent colonists in human gut metatranscriptomes, then a petabase-scale
240 SRA search expanded this to 2,152 species genomes (defined here as Oblin-1 90% aaid clusters) [9]. With
241 Logan we detected an additional 2,964 Obelisk species, a 2.4-fold increase. Likewise, the total Obelisk
242 count increased 3.7-fold, to 26,263 sequences (Fig. 4b). The expanded Oblin-1 proteins have no homologs
243 in NCBI nr, and no homologs known outside of Obelisks. This expanded dataset should accelerate research
244 to uncover the role of these mysterious elements.

245 **P4 Satellites** We then searched for P4-like satellites [34], mobile elements that hijack bacteriophages.
246 P4-like satellites were recently found to be numerous in enterobacterial genomes, where they provide
247 the host with anti-phage functions. We observed a 4.5-fold expansion of the number of these satellites
248 (Fig. 4c), including elements that are unrelated with previous sub-families (Extended Data Fig. 8c). The
249 pan-genome of these elements is thus doubled in relation to RefSeq (Extended Data Fig. 8c). These newly
250 uncovered P4 elements may thus encode many novel anti-phage functions of ecological and biotechnological
251 relevance. These searches illustrate how to leverage Logan to mine for complex genetic elements with
252 multiple core genes.

253 **Plasmids** Plasmids are mobile genetic elements of Bacteria and Archaea that play a critical role in
254 horizontal gene transfer, driving processes such as the spread of antibiotic resistance genes. Given their
255 ecological and clinical importance, systematically characterizing plasmid diversity can provide valuable
256 insights into global gene transfer patterns. We extracted all circular contigs from Logan assemblies and
257 identified 468,614 putatively complete plasmids (264,160 unique sequences) across 195,347 metagenome,
258 57,148 bacterial, and 12 archaeal isolate accessions. To evaluate the global distribution of these plasmids,
259 including in samples where they were not fully assembled as circular contigs, we first reduced redundancy
260 in the dataset by grouping highly similar sequences into 60,331 clusters, and then mapped representative
261 sequences from each cluster to all Logan contigs. This approach identified plasmids in 2,095,914 samples
262 over the globe (Fig. 4d), highlighting their widespread distribution. Moreover, the number of distinct
263 detectable plasmids has steadily increased over time (Fig. 4e).

264 Next, we investigated the origins underlying this extensive plasmid diversity by analyzing the compo-
265 sition of plasmid clusters. We found that 92.5% comprised only metagenomic sequences, while just 3.6%
266 consisted exclusively of plasmids from cultured organisms. Despite this predominance of environmental
267 sequences, only 17.5% of metagenomic plasmids could be assigned to known replicon families compared
268 to 78.3% from isolates, highlighting that the plasmid diversity in natural environments remains largely
269 uncharacterized. Environmental plasmids were depleted of known antimicrobial resistance (AMR) genes
270 and encoded more antimicrobial peptides (AMPs) relative to plasmids from isolates (Fig. 4f), reflecting
271 distinct accessory gene repertoires and underscoring the value of assessing plasmid diversity across diverse
272 sample types, as enabled by Logan.

273 To quantify the extent of plasmid diversity uncovered by surveying Logan assemblies, we measured
274 the phylogenetic diversity of selected replicase and relaxase proteins from these plasmids and compared
275 it with complete plasmid genomes from established databases [35, 36]. Plasmids identified from Logan
276 assemblies expanded the phylogenetic diversity up to 21.8-fold compared to PLSDB (Fig. 4g) and up to
277 7.0-fold compared to IMG/PR (Data Availability, Plasmid PD Table). Overall, these results demonstrate
278 that mining Logan assemblies reveals a vast and previously undiscovered diversity of genetic elements not
279 captured by other genomic data resources.

280 **Antimicrobial resistance** Logan also enables the global-scale analysis of antimicrobial resistance
281 (AMR) across all publicly available sequencing data. By aligning all Logan contigs to the CARD database,

we identified 7.9 million AMR-positive (AMR+) SRA accessions and 13,000 AMR+ plasmids (Extended Data Fig. 7a). AMR genes are enriched in metagenomes across the SRA, whereas plasmids show the opposite pattern, driven by the large fraction of bacterial isolates in plasmid datasets (Extended Data Fig. 7b). We observe an enrichment of human and livestock metagenomic samples in AMR-positive datasets (Extended Data Fig. 7c), both in SRA accessions and plasmids. Geographically, AMR+ metagenomes are broadly distributed, and their discovery has increased steadily over the past two decades based on collection dates (Extended Data Fig. 7d-e). AMR gene content varies across metagenome categories, with wastewater, livestock, and human metagenomes showing the highest enrichment in AMR gene counts per accession (Extended Data Fig. 7f-g). As sequencing databases continue to grow, full-scale sequence indexes can be re-purposed as an AMR-surveillance network.

2 Discussion

The exponential growth of sequence databases has created a paradox: humanity is generating more genomics data than ever before, which is making the data increasingly inaccessible due to computational barriers. Logan resolves this paradox through transforming raw sequencing data into accessible and searchable resources which enable systematic analysis of global sequence diversity, and offers a technical framework by which analysis capacity can continue to scale alongside database growth.

Earth's genetic diversity is a heritage of humanity [37]. To bulwark against a trend of commercializing public scientific data, we release all Logan data into the public-domain, and emphasize the continued need for community development of free and unrestricted data commons (<https://registry.opendata.aws/pasteur-logan/>). For decades, BLAST democratized sequence comparison at the gigabase scale, and it transformed biology research. Logan brings such a capacity to the petabase-era, enabling analogous discoveries across orders of magnitude larger datasets. Like the original NCBI web server that made BLAST a ubiquitous tool, Logan-Search's web-interface (<https://logan-search.org/>) ensures that this resource is practically available for the research community.

Logan-Search enables researchers to rapidly test hypotheses across the breadth of public sequencing data, uncovering unexpected connections that were previously hidden in plain sight, with direct applications to biotechnology and planetary health. The discovery of HHV-6 reactivation in therapeutic CAR T-cells illustrated how large-scale sequence search could inform clinically relevant insights [22]. Here, we generalize such capability and extend this discovery to characterize latent viral reactivation tumor-infiltrating lymphocytes. Logan enables researchers to query the sampled biosphere for specific functions, such as identifying novel plastic-degrading enzymes that outperform engineered variants. It allows us to directly sample from nature's vast parallel evolutionary experiment.

The massive expansion of plastic-active enzyme homologs makes it immediately obvious that Logan enables the free and unprecedented capacity to screen for variants and novel versions of proteins and genes for biotechnology. This includes but is not limited to, for example, identifying novel viral vectors, or receptors/effectors with altered tropism; biosynthetic gene clusters for the production of antibiotics or natural products; efficient or process-optimized industrial enzymes such as proteases, amylases, or cellulases; or biotechnology enzymes such as Cas, reverse-transcriptases, or polymerases. Moreover, these public datasets are a rich resource with obvious applications as training data for a next-generation of machine learning and artificial intelligence models.

The ability to efficiently analyze all public sequence data arrives at a critical moment for biodiversity research. Current sampling of Earth's genetic diversity shows clear geographic and taxonomic biases, evident in SRA metadata. As climate change and habitat loss accelerate species extinction, systematic sequence analysis becomes essential not only for documenting disappearing diversity but for understanding the genetic basis of adaptation and resilience. Logan provides the technical foundations, while highlighting the urgent need for broader and more representative sampling of Earth's biosphere.

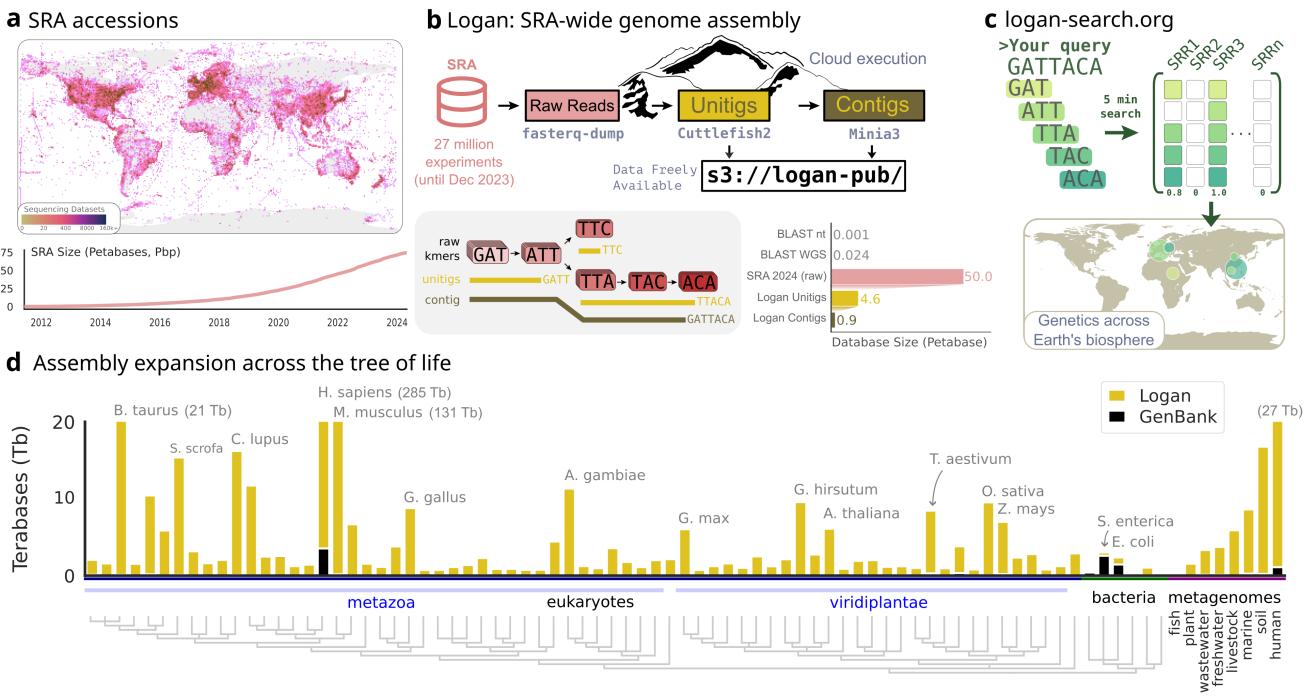


Figure 1: Assembling all accessions of the SRA using a cloud architecture into unitigs and contigs. (a) Geographic distribution of samples over the Sequence Read Archive (SRA), and the near-exponential growth of SRA in terms of number of cumulative accession size of raw data. (b) Top diagram describes the cloud computation workflow of Logan, starting from SRA reads, then computing unitigs and contigs assemblies, and finally uploading data to our public repository. Bottom left diagram shows a toy dataset with k -mers extracted from raw reads, then unitigs and contigs constructed. Bottom right bar plot represents the size of the SRA compared to Logan assembled unitigs and contigs in sum of bases, and WGS and BLAST databases. (c) The logan-search.org service enables searching an arbitrary query (example: “GATTACA”) against the full unitig index of the SRA in less than 5 min; hits are mapped to their geographic origins. (d) Tree of Life sampled with the 116 most abundant taxa from NCBI GenBank WGS as well as 116 most abundant taxa in Logan assemblies, according to NCBI taxonomy. Black bars represent the total number of assembled bases in GenBank WGS, and yellow bars the additional number of bases in Logan contigs. Bars exceeding 20 terabases are capped and their true total assembly size is annotated. Assembled bases for a subset of metagenome types are represented separately as the 8 rightmost bars.

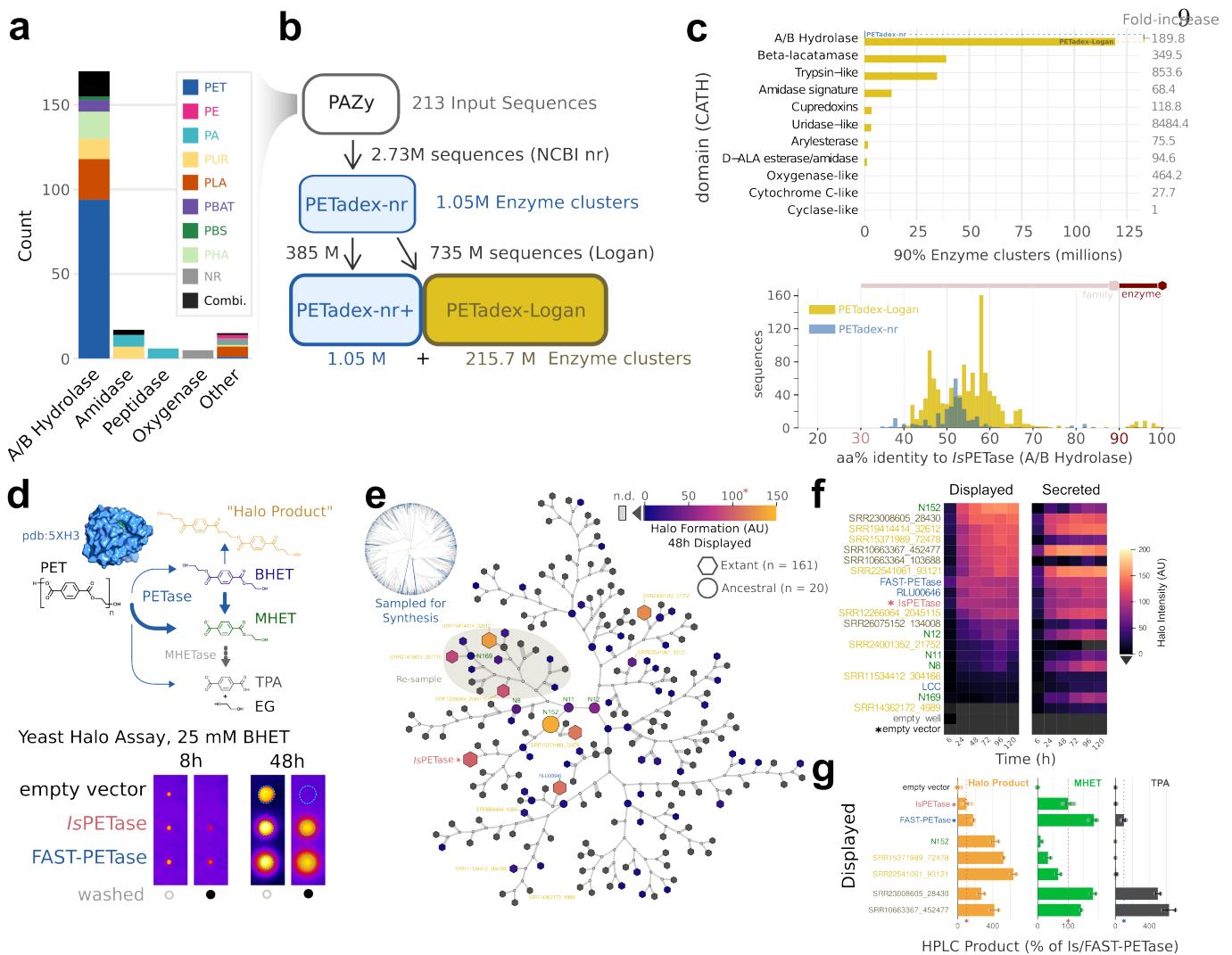


Figure 2: Discovering novel and efficacious plastic-active enzymes. (a) The domains and activity of the 213 experimentally validated plastic-active enzyme (PAZy) search query (Extended Data Fig. 2a). (b) Logan *PETadex* homology search returned 216.75 million PAZy-homologs after clustering at 90% amino acid identity (Extended Data Fig. 2b,c). (c) *PETadex-Logan* is a >200-fold expansion of candidate PAZy relative to NCBI nr across distinct PAZy CATH domains (see Methods). Histogram shows the distribution of *IsPETase*-aligned sequences, illustrating that Logan (yellow) uncovers more diversity across the detectable range of sequence identities relative to NCBI nr (blue). (d) The PETase reaction which underpins the high-throughput yeast-based halo assay. Yeast expressing either control (*IsPETase*) or candidate enzyme targeting the PET substrate BHET were grown on agar plates to create a white halo which is quantified as pixel intensity (shown as pseudocolored), before (open circle) and after washing (black circle) around the colony (cyan outline) (Extended Data Fig. 3a,b). High-performance liquid chromatography (HPLC) and mass spectroscopy suggest that the “halo product” is O,O’-(ethane-1,2-diy) bis(oxy(2-hydroxyethyl)carbonyl)terephthalate (Extended Data Fig. 3e,f). (e) Phylogenetic tree of sampled candidate PAZy that were synthesized and experimentally screened. Nodes are colored based on 48 hour halo formation activity in surface-displayed expression. The gray-highlighted clade was re-sampled for additional sequences. (f) Heatmap of select enzyme halo formation activity over time, quantified in surface-display and secreted systems (Extended Data Fig. 4a). (g) Quantitative validation of candidate high-activity *PETadex-Logan* enzymes by HPLC. The bars show the percentage of product formed relative to the activity of *IsPETase* (halo product, MHET) or FAST-PETase (TPA). Logan enzymes demonstrate product formation exceeding that of *IsPETase* and FAST-PETase.

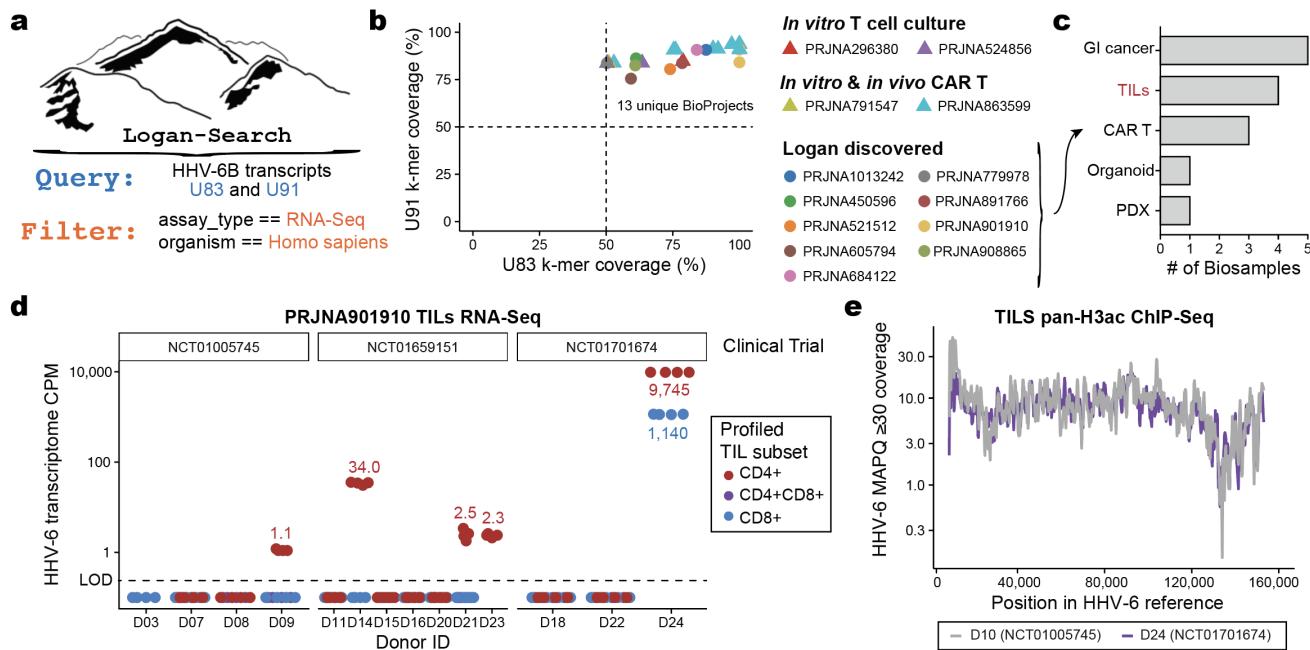


Figure 3: Identification and reactivation of HHV-6 in large-scale RNA-seq datasets. (a) Input query to Logan-Search for two abundant HHV-6 genes (U83 and U91) and filtering criteria for human RNA-seq BioSamples. (b) K-mer coverage analysis of 13 identified HHV-6-positive BioProjects, including 9 with no prior HHV-6 annotation (circles). Triangles indicating previously annotated datasets (Serratus and HHV6 paper). (c) Annotation of newly discovered HHV-6-bearing BioSamples, including gastrointestinal cancers, tumor-infiltrating lymphocytes (TILs), chimeric antigen receptor (CAR) T-cell products, organoids, and patient-derived xenografts (PDXs). (d) RNA-seq analyses of TIL cultures (PRJNA901910). Values indicate HHV-6 RNA abundance (counts per million, CPM) out of the full library, reflecting HHV-6 reactivation from cultured T cells. (e) ChIP-seq analyses of TIL cultures (PRJNA901909). Values reflect the HHV-6 DNA abundance for two donors in H3 acetylation chromatin, reflecting coverage across the full viral contig.

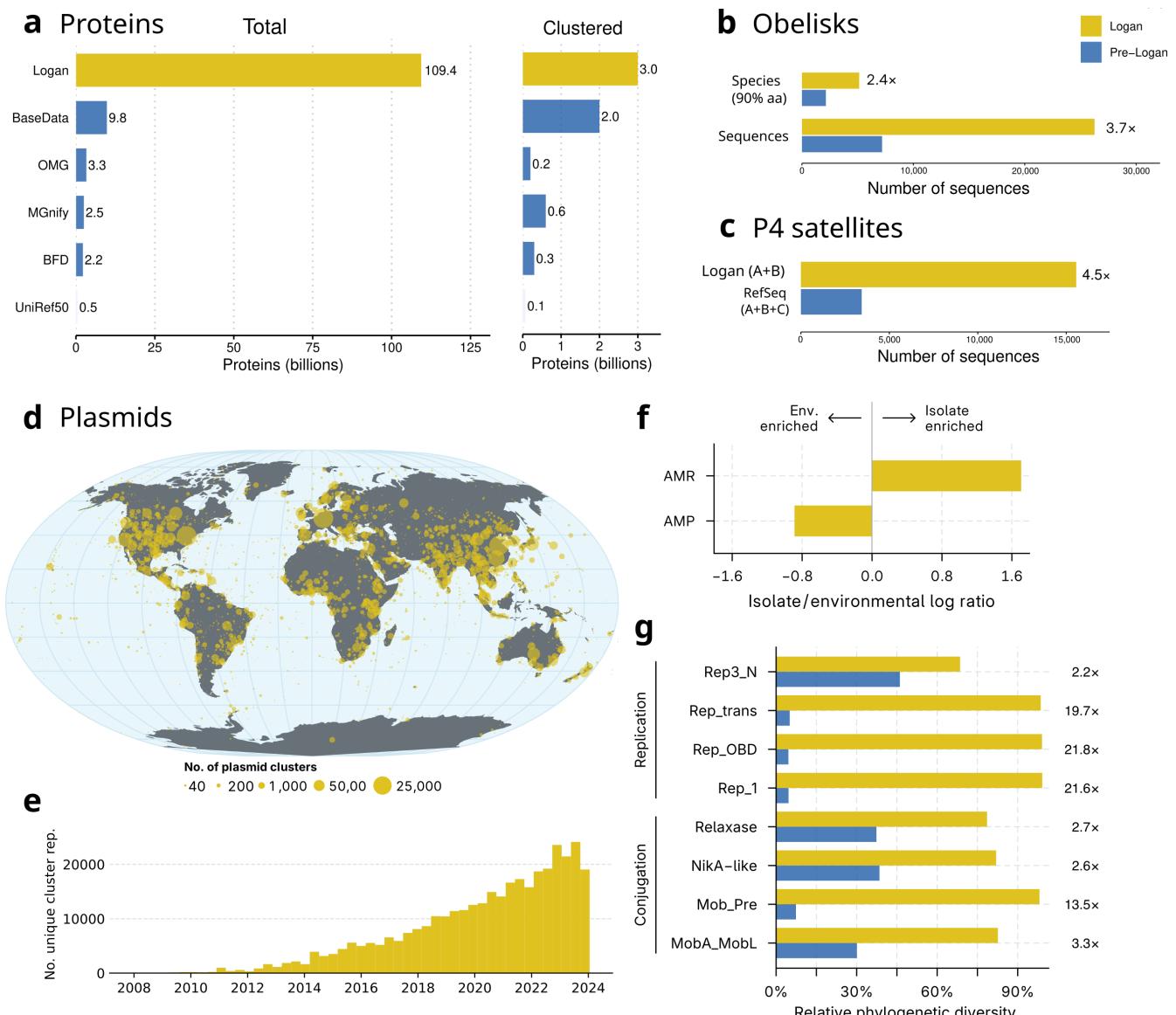


Figure 4: Expanding the Known Universe of Proteins, Plasmids, and Viral Elements (a) Bar plot labeled “Total” shows the total number of proteins extracted from Logan contigs, compared to other databases, and bar plot labeled “Clustered” shows the same set of proteins but clustered at 50% identity. (b) Expansion of Obelisks requiring circular contigs with full-length Oblin-1 proteins, identified in Logan (yellow), relative to the initial petabase-scale search [9] (blue). Total sequences and species (clustered centroids of Oblin-1) are shown. (c) Number of P4 satellites found in Logan contigs (types A+B) compared to those in RefSeq (types A+B+C). Types A, B, C refer to the number of core components detected by SatelliteFinder (Methods): A = all 7, B = 6 out of 7, C = 5 out of 7. (d) Geographical distribution of plasmids detected over 2,095,914 samples with geolocation data. Circle areas are proportional to the number of distinct plasmid clusters per region. For visual clarity, spatially close samples were grouped using DBSCAN, and circles are placed at the coordinates of the corresponding cluster medoids. The map uses the Loximuthal projection. (e) Number of distinct plasmid cluster representatives detected over time. Counts (y-axis) are shown in 120-day intervals (x-axis). (f) Comparison of accessory gene repertoires in plasmids from environmental samples vs. cultured isolates. Plasmids from isolates are comparatively enriched for antimicrobial resistance (AMR) genes, whereas plasmids from environmental sources are enriched for antimicrobial peptides (AMPs). Functional enrichment (x-axis) was quantified as the ratio of gene density (genes per megabase) for a given function (AMR or AMP) between the two plasmid groups. (g) Relative Faith’s phylogenetic diversity of selected replicases and relaxases encoded by plasmids identified in Logan (yellow) and those retrieved from PLSDB (blue). The phylogenetic diversity fold change, indicated by numbers to the right of the bars, represents the ratio of the diversity across all plasmids (Logan and PLSDB combined) to the diversity in PLSDB plasmids alone.

328 References

- 329 [1] Kenneth Katz, Oleg Shutov, Richard Lapoint, Michael Kimelman, J Rodney Brister, and Christopher
330 O’Sullivan. The Sequence Read Archive: a decade more of explosive growth. *Nucleic Acids Research*,
331 50(D1):D387–D390, 2022.
- 332 [2] Robert C Edgar, Brie Taylor, Victor Lin, Tomer Altman, Pierre Barbera, Dmitry Meleshko, Dan
333 Lohr, Gherman Novakovsky, Benjamin Buchfink, Basem Al-Shayeb, et al. Petabase-scale sequence
334 alignment catalyses viral discovery. *Nature*, 602(7895):142–147, 2022.
- 335 [3] Phelim Bradley, Henk C Den Bakker, Eduardo PC Rocha, Gil McVean, and Zamin Iqbal. Ultrafast
336 search of all deposited bacterial and viral genomic data. *Nature Biotechnology*, 37(2):152–159, 2019.
- 337 [4] Mikhail Karasikov, Harun Mustafa, Daniel Danciu, Christopher Barber, Marc Zimmermann, Gunnar
338 Rätsch, and André Kahles. Metagraph: Indexing and analysing nucleotide archives at petabase-scale.
339 *BioRxiv*, pages 2020–10, 2020.
- 340 [5] Kenneth S Katz, Oleg Shutov, Richard Lapoint, Michael Kimelman, J Rodney Brister, and Christopher
341 O’Sullivan. STAT: a fast, scalable, MinHash-based k-mer tool to assess Sequence Read Archive
342 next-generation sequence submissions. *Genome Biology*, 22:1–15, 2021.
- 343 [6] Sergey A Shiryev and Richa Agarwala. Indexing and searching petabase-scale nucleotide resources.
344 *Nature Methods*, pages 1–9, 2024.
- 345 [7] Martin Hunt, Leandro Lima, Wei Shen, John Lees, and Zamin Iqbal. AllTheBacteria-all bacterial
346 genomes assembled, available and searchable. *bioRxiv*, pages 2024–03, 2024.
- 347 [8] Luiz Irber, N Tessa Pierce-Ward, and C Titus Brown. Sourmash branchwater enables lightweight
348 petabyte-scale sequence search. *bioRxiv*, pages 2022–11, 2022.
- 349 [9] Ivan N Zheludev, Robert C Edgar, Maria Jose Lopez-Galiano, Marcos De la Pena, Artem Babaian,
350 Ami S Bhatt, and Andrew Z Fire. Viroid-like colonists of human microbiomes. *Cell*, 187(23):6521–
351 6536, 2024.
- 352 [10] Eric W Sayers, Jeff Beck, Evan E Bolton, J Rodney Brister, Jessica Chan, Donald C Comeau, Ryan
353 Connor, Michael DiCuccio, Catherine M Farrell, Michael Feldgarden, et al. Database resources of
354 the national center for biotechnology information. *Nucleic Acids Research*, 52(D1):D33, 2024.
- 355 [11] Benjamin Buchfink, Klaus Reuter, and Hajk-Georg Drost. Sensitive protein alignments at tree-of-life
356 scale using DIAMOND. *Nature Methods*, 18(4):366–368, 2021.
- 357 [12] Asher Mullard. FDA approves first tumour-infiltrating lymphocyte (TIL) therapy, bolstering hopes
358 for cell therapies in solid cancers. *Nat Rev Drug Discov*, 23(4):238, 2024.
- 359 [13] Lorna Richardson, Ben Allen, Germana Baldi, Martin Beracochea, Maxwell L Bileschi, Tony Bur-
360 dett, Josephine Burgin, Juan Caballero-Pérez, Guy Cochrane, Lucy J Colwell, et al. MGnify: the
361 microbiome sequence data analysis resource in 2023. *Nucleic Acids Research*, 51(D1):D753–D759,
362 2023.
- 363 [14] Richard C Thompson, Winnie Courtene-Jones, Julien Boucher, Sabine Pahl, Karen Raubenheimer,
364 and Albert A Koelmans. Twenty years of microplastic pollution research—what have we learned?
365 *Science*, 386(6720):eadl2746, 2024.

- 366 [15] Philip J Landrigan, Sarah Dunlop, Marina Treskova, Hervé Raps, Christos Symeonides, Jane Muncke,
367 Margaret Spring, John Stegeman, Bethanie Carney Almroth, Thomas C Chiles, Maureen Cropper,
368 Megan Deeney, Lizzie Fuller, Roland Geyer, Rachel Karasik, Tiza Mafira, Alexander Mangwiwo,
369 Denise Margaret Matias, Yannick Mulders, Yongjoon Park, Costas A Velis, Roel Vermeulen, Martin
370 Wagner, Zhanyung Wang, Ella M Whitman, Tracey J Woodruff, and Joacim Rocklöv. *The Lancet*, 2025.
371
- 372 [16] Shosuke Yoshida, Kazumi Hiraga, Toshihiko Takehana, Ikuo Taniguchi, Hironao Yamaji, Yasuhito
373 Maeda, Kiyotsuna Toyohara, Kenji Miyamoto, Yoshiharu Kimura, and Kohei Oda. A bacterium that
374 degrades and assimilates poly (ethylene terephthalate). *Science*, 351(6278):1196–1199, 2016.
- 375 [17] Hongyuan Lu, Daniel J Diaz, Natalie J Czarnecki, Congzhi Zhu, Wantae Kim, Raghav Shroff,
376 Daniel J Acosta, Bradley R Alexander, Hannah O Cole, Yan Zhang, et al. Machine learning-aided
377 engineering of hydrolases for pet depolymerization. *Nature*, 604(7907):662–667, 2022.
- 378 [18] Patrick CF Buchholz, Golo Feuerriegel, Hongli Zhang, Pablo Perez-Garcia, Lena-Luisa Nover, Jen-
379 nifer Chow, Wolfgang R Streit, and Jürgen Pleiss. Plastics degradation by hydrolytic enzymes:
380 The plastics-active enzymes database—PAZy. *Proteins: Structure, Function, and Bioinformatics*,
381 90(7):1443–1456, 2022.
- 382 [19] Jianwei Chen, Yangyang Jia, Ying Sun, Kun Liu, Changhao Zhou, Chuan Liu, Denghui Li, Guilin
383 Liu, Chengsong Zhang, Tao Yang, et al. Global marine microbial diversity and its potential in
384 bioprospecting. *Nature*, 633(8029):371–379, 2024.
- 385 [20] Hogyun Seo, Hwaseok Hong, Jiyoung Park, Seul Hoo Lee, Dongwoo Ki, Aejin Ryu, Hye-Young
386 Sagong, and Kyung-Jin Kim. Landscape profiling of PET depolymerases using a natural sequence
387 cluster framework. *Science*, 387(6729):eadp5637, 2025.
- 388 [21] Ian Sillitoe, Nicola Bordin, Natalie Dawson, Vaishali P Waman, Paul Ashford, Harry M Scholes,
389 Camilla SM Pang, Laurel Woodridge, Clemens Rauer, Neeladri Sen, et al. CATH: increased structural
390 coverage of functional space. *Nucleic Acids Research*, 49(D1):D266–D273, 2021.
- 391 [22] Caleb A Lareau, Yajie Yin, Katie Maurer, Katalin D Sandor, Bence Daniel, Garima Yagnik, José
392 Peña, Jeremy Chase Crawford, Anne M Spanjaart, Jacob C Gutierrez, et al. Latent human her-
393 pesvirus 6 is reactivated in CAR T cells. *Nature*, 623(7987):608–615, 2023.
- 394 [23] Hiu Tung Chow, Wenjing Li, Bin Yang, Friedrich von Wintzingerode, and Qi Chen. HHV-6 and
395 HHV-7 reactivation in allogeneic CAR-T cell therapy. *Trends in Biotechnology*, 2025.
- 396 [24] Sarah A LaMere, Ryan C Thompson, H Kiyomi Komori, Adam Mark, and Daniel R Salomon.
397 Promoter H3K4 methylation dynamically reinforces activation-induced pathways in human CD4 T
398 cells. *Genes & Immunity*, 17(5):283–297, 2016.
- 399 [25] Iart Luca Shytaj, Bojana Lucic, Mattia Forcato, Carlotta Penzo, James Billingsley, Vibor Laketa,
400 Steven Bosinger, Mia Stanic, Francesco Gregoretti, Laura Antonelli, et al. Alterations of redox and
401 iron metabolism accompany the development of HIV latency. *The EMBO journal*, 39(9):e102209,
402 2020.
- 403 [26] Yongxian Hu, Yali Zhou, Mingming Zhang, Houli Zhao, Guoqing Wei, Wengang Ge, Qu Cui, Qitian
404 Mu, Gong Chen, Lu Han, et al. Genetically modified CD7-targeting allogeneic CAR-T cell therapy
405 with enhanced efficacy for relapsed/refractory CD7-positive hematological malignancies: a phase I
406 clinical study. *Cell research*, 32(11):995–1007, 2022.

- 407 [27] Eva Eliassen, Emily Lum, Joshua Pritchett, Joseph Ongradi, Gerhard Krueger, John R Crawford,
408 Tuan L Phan, Dharam Ablashi, and Stanley David Hudnall. Human herpesvirus 6 and malignancy:
409 a review. *Frontiers in oncology*, 8:512, 2018.
- 410 [28] Zongcheng Li, Lei Zhao, Yuanyuan Zhang, Li Zhu, Wei Mu, Tong Ge, Jin Jin, Jiaqi Tan, Jiali Cheng,
411 Jue Wang, et al. Functional diversification and dynamics of CAR-T cells in patients with B-ALL.
412 *Cell Reports*, 42(10), 2023.
- 413 [29] Shannon S Choi, Vincent van Unen, Huimin Zhang, Arjun Rustagi, Samira A Alwahabi, António JM
414 Santos, Joshua E Chan, Brandon Lam, Daniel Solis, Jordan Mah, et al. Organoid modeling of lung-
415 resident immune responses to SARS-CoV-2 infection. *Research Square*, pages rs-3, 2023.
- 416 [30] Brian Thompson, Ann Strange, Carol M Amato, Jonathan Hester-McCullough, Amod A Sarnaik,
417 Jeffrey S Weber, and David M Woods. CD4 phenotypes are associated with reduced expansion of
418 tumor-infiltrating lymphocytes in melanoma patients treated with adoptive cell therapy. *The Journal
419 of Immunology*, 211(5):735–742, 2023.
- 420 [31] Oliver Vince, Phoebe Oldach, Valerio Pereno, Marcus HY Leung, Carla Greco, Gus Minto-Cowcher,
421 Saif Ur-Rehman, Keith YK Kam, William Chow, Emma Bolton, et al. Breaking through biology's
422 data wall: Expanding the known tree of life by over 10x using a global biodiscovery pipeline. *bioRxiv*,
423 pages 2025–06, 2025.
- 424 [32] Andre Cornman, Jacob West-Roberts, Antonio Pedro Camargo, Simon Roux, Martin Beracochea,
425 Milot Mirdita, Sergey Ovchinnikov, and Yunha Hwang. The OMG dataset: An Open MetaGenomic
426 corpus for mixed-modality genomic language modeling. *bioRxiv*, pages 2024–08, 2024.
- 427 [33] Milot Mirdita, Konstantin Schütze, Yoshitaka Moriwaki, Lim Heo, Sergey Ovchinnikov, and Martin
428 Steinegger. ColabFold: making protein folding accessible to all. *Nature Methods*, 19(6):679–682,
429 2022.
- 430 [34] Jorge A Moura de Sousa, Alfred Fillol-Salom, José R Penadés, and Eduardo P C Rocha. Identification
431 and characterization of thousands of bacteriophage satellites across bacteria. *Nucleic Acids Research*,
432 51(6):2759–2777, April 2023.
- 433 [35] Leidy-Alejandra G Molano, Pascal Hirsch, Matthias Hannig, Rolf Müller, and Andreas Keller. The
434 PLSDB 2025 update: enhanced annotations and improved functionality for comprehensive plasmid
435 research. *Nucleic Acids Research*, 53(D1):D189–D196, January 2025.
- 436 [36] Antonio Pedro Camargo, Lee Call, Simon Roux, Stephen Nayfach, Marcel Huntemann, Krishnaveni
437 Palaniappan, Anna Ratner, Ken Chu, Supratim Mukherjee, T B K Reddy, I-Min A Chen, Natalia N
438 Ivanova, Emiley A Elo-Fadrosh, Tanja Woyke, David A Baltrus, Salvador Castañeda-Barba, Fer-
439 nando de la Cruz, Barbara E Funnell, James P J Hall, Aindrila Mukhopadhyay, Eduardo P C Rocha,
440 Thibault Stalder, Eva Top, and Nikos C Kyripies. IMG/PR: a database of plasmids from genomes
441 and metagenomes with rich annotations and metadata. *Nucleic Acids Research*, 52(D1):D164–D173,
442 January 2024.
- 443 [37] Soren Brunak, Antoine Danchin, Masahira Hattori, Haruki Nakamura, Kazuo Shinozaki, Tara Matise,
444 and Daphne Preuss. Nucleotide sequence database policies. *Science*, 298(5597):1333–1333, 2002.
- 445 [38] Jamshed Khan, Marek Kokot, Sebastian Deorowicz, and Rob Patro. Scalable, ultra-fast, and low-
446 memory construction of compacted de Bruijn graphs with Cuttlefish 2. *Genome Biology*, 23(1):190,
447 2022.
- 448 [39] Marek Kokot, Maciej Dlugosz, and Sebastian Deorowicz. KMC 3: counting and manipulating k-mer
449 statistics. *Bioinformatics*, 33(17):2759–2761, 2017.

- 450 [40] Rayan Chikhi and Guillaume Rizk. Space-efficient and exact de Bruijn graph representation based
451 on a Bloom filter. *Algorithms for Molecular Biology*, 8:1–9, 2013.
- 452 [41] Andrey Prjibelski, Dmitry Antipov, Dmitry Meleshko, Alla Lapidus, and Anton Korobeynikov. Using
453 SPAdes de novo assembler. *Current Protocols in Bioinformatics*, 70(1), June 2020.
- 454 [42] Cécile Monat, Sudharsan Padmarasu, Thomas Lux, Thomas Wicker, Heidrun Gundlach, Axel Him-
455 melbach, Jennifer Ens, Chengdao Li, Gary J Muehlbauer, Alan H Schulman, et al. TRITEX:
456 chromosome-scale sequence assembly of Triticeae genomes with open-source tools. *Genome Biol-*
457 *ogy*, 20(1):284, 2019.
- 458 [43] Alex Di Genova, Elena Buena-Atienza, Stephan Ossowski, and Marie-France Sagot. Efficient hybrid
459 de novo assembly of human genomes with WENGAN. *Nature Biotechnology*, 39(4):422–430, 2021.
- 460 [44] Thomas Krannich, W Timothy J White, Sebastian Niehus, Guillaume Holley, Bjarni V Halldórsson,
461 and Birte Kehr. Population-scale detection of non-reference sequence variants using colored de Bruijn
462 graphs. *Bioinformatics*, 38(3):604–611, 2022.
- 463 [45] Rayan Chikhi, Antoine Limasset, and Paul Medvedev. Compacting de Bruijn graphs from sequencing
464 data quickly and in low memory. *Bioinformatics*, 32(12):i201–i208, 2016.
- 465 [46] Annika Jochheim, Florian E Jochheim, Alexandra Kolodyazhnaya, Etienne Morice, Martin Steineg-
466 ger, and Johannes Soeding. Strain-resolved de-novo metagenomic assembly of viral genomes and
467 microbial 16S rRNAs. *bioRxiv*, pages 2024–03, 2024.
- 468 [47] Dmitry Meleshko, Iman Hajirasouliha, and Anton Korobeynikov. coronaSPAdes: from biosynthetic
469 gene clusters to RNA viral assemblies. *Bioinformatics*, 38(1):1–8, 2022.
- 470 [48] Yann Collet. Rfc 8878: Zstandard compression and the 'application/zstd' media type, 2021.
- 471 [49] Wei Shen, Shuai Le, Yan Li, and Fuquan Hu. SeqKit: a cross-platform and ultrafast toolkit for
472 FASTA/Q file manipulation. *PloS One*, 11(10):e0163962, 2016.
- 473 [50] Heng Li. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*, 34(18):3094–3100,
474 September 2018.
- 475 [51] Téo Lemane, Paul Medvedev, Rayan Chikhi, and Pierre Peterlongo. Kmtricks: efficient and flex-
476 ible construction of Bloom filters for large sequencing data collections. *Bioinformatics Advances*,
477 2(1):vbac029, 2022.
- 478 [52] Téo Lemane, Nolan Lezzoche, Julien Lecubin, Eric Pelletier, Magali Lescot, Rayan Chikhi, and Pierre
479 Peterlongo. Indexing and real-time user-friendly queries in terabyte-sized complex genomic datasets
480 with kmindex and ORA. *Nature Computational Science*, 4(2):104–109, 2024.
- 481 [53] Burton H Bloom. Space/time trade-offs in hash coding with allowable errors. *Communications of
482 the ACM*, 13(7):422–426, 1970.
- 483 [54] Lucas Robidou and Pierre Peterlongo. findere: fast and precise approximate membership query.
484 In *String Processing and Information Retrieval: 28th International Symposium, SPIRE 2021, Lille,
485 France, October 4–6, 2021, Proceedings 28*, pages 151–163. Springer, 2021.
- 486 [55] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Os-
487 trow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint
488 arXiv:2410.21276*, 2024.

- 489 [56] Anthony Baire, Pierre Marijon, Francesco Andreace, and Pierre Peterlongo. Back to sequences: Find
490 the origin of k-mers. *Journal of Open Source Software*, 9(101):7066, 2024.
- 491 [57] Ying Chen, Weicai Ye, Yongdong Zhang, and Yuesheng Xu. High speed BLASTN: an accelerated
492 MegaBLAST search tool. *Nucleic Acids Research*, 43(16):7762–7768, August 2015.
- 493 [58] Isuru Karunatillaka, Lukasz Jaroszewski, and Adam Godzik. Novel putative polyethylene tereph-
494 thalate (PET) plastic degrading enzymes from the environmental metagenome. *Proteins: Structure,*
495 *Function, and Bioinformatics*, 90(2):504–511, 2022.
- 496 [59] Michel van Kempen, Stephanie S. Kim, Charlotte Tumescheit, Milot Mirdita, Jeongjae Lee, Cameron
497 L. M. Gilchrist, Johannes Söding, and Martin Steinegger. Fast and accurate protein structure search
498 with Foldseek. *Nature Biotechnology*, 42(2):243–246, February 2024.
- 499 [60] Robert C. Edgar. Search and clustering orders of magnitude faster than BLAST. *Bioinformatics*,
500 26(19):2460–2461, October 2010.
- 501 [61] Josh Abramson, Jonas Adler, Jack Dunger, Richard Evans, Tim Green, Alexander Pritzel, Olaf
502 Ronneberger, Lindsay Willmore, Andrew J. Ballard, Joshua Bambrick, Sebastian W. Bodenstein,
503 David A. Evans, Chia-Chun Hung, Michael O'Neill, David Reiman, Kathryn Tunyasuvunakool,
504 Zachary Wu, Akvilė Žemgulytė, Eirini Arvaniti, Charles Beattie, Ottavia Bertolli, Alex Bridgland,
505 Alexey Cherepanov, Miles Congreve, Alexander I. Cowen-Rivers, Andrew Cowie, Michael Figurnov,
506 Fabian B. Fuchs, Hannah Gladman, Rishabh Jain, Yousuf A. Khan, Caroline M. R. Low, Kuba Perlin,
507 Anna Potapenko, Pascal Savy, Sukhdeep Singh, Adrian Stecula, Ashok Thillaisundaram, Catherine
508 Tong, Sergei Yakneen, Ellen D. Zhong, Michal Zielinski, Augustin Žídek, Victor Bapst, Pushmeet
509 Kohli, Max Jaderberg, Demis Hassabis, and John M. Jumper. Accurate structure prediction of
510 biomolecular interactions with alphafold 3. *Nature*, 630(8016):493–500, May 2024.
- 511 [62] Robert C. Edgar. Muscle5: High-accuracy alignment ensembles enable unbiased assessments of
512 sequence homology and phylogeny. *Nature Communications*, 13(1):6968, November 2022.
- 513 [63] Sean R. Eddy. Accelerated Profile HMM Searches. *PLoS Computational Biology*, 7(10):e1002195,
514 October 2011.
- 515 [64] Bui Quang Minh, Heiko A Schmidt, Olga Chernomor, Dominik Schrempf, Michael D Woodhams,
516 Arndt von Haeseler, and Robert Lanfear. IQ-tree 2: New models and efficient methods for phylo-
517 genetic inference in the genomic era. *Molecular Biology and Evolution*, 37(5):1530–1534, February
518 2020.
- 519 [65] Gabriel Foley, Ariane Mora, Connie M Ross, Scott Bottoms, Leander Sütl, Marnie L Lamprecht,
520 Julian Zaugg, Alexandra Essebier, Brad Balderson, Rhys Newell, et al. Engineering indel and substi-
521 tution variants of diverse and ancient enzymes using Graphical Representation of Ancestral Sequence
522 Predictions (GRASP). *PLoS computational biology*, 18(10):e1010633, 2022.
- 523 [66] Carrie Baker Brachmann, Adrian Davies, Gregory J Cost, Emerita Caputo, Joachim Li, Philip
524 Hieter, and Jef D Boeke. Designer deletion strains derived from *Saccharomyces cerevisiae* S288C: a
525 useful set of strains and plasmids for PCR-mediated gene disruption and other applications. *Yeast*,
526 14(2):115–132, 1998.
- 527 [67] Colin JB Harvey, Mancheng Tang, Ulrich Schlecht, Joe Horecka, Curt R Fischer, Hsiao-Ching Lin,
528 Jian Li, Brian Naughton, James A Cherry, Molly Miranda, et al. Hex: A heterologous expression
529 platform for the discovery of fungal natural products. *Science advances*, 4(4):eaar5459, 2018.

- 530 [68] Leslie A Mitchell, James Chuang, Neta Agmon, Chachrit Khunsriraksakul, Nick A Phillips, Yizhi
531 Cai, David M Truong, Ashan Veerakumar, Yuxuan Wang, Maria Mayorga, et al. Versatile genetic as-
532 sembly system (VEGAS) to assemble pathways for expression in *S. cerevisiae*. *Nucleic acids research*,
533 43(13):6620–6630, 2015.
- 534 [69] Valentina Pirillo, Loredano Pollegioni, and Gianluca Molla. Analytical methods for the investigation
535 of enzyme-catalyzed degradation of polyethylene terephthalate. *The FEBS Journal*, 288(16):4730–
536 4745, 2021.
- 537 [70] Nicolas L Bray, Harold Pimentel, Pál Melsted, and Lior Pachter. Near-optimal probabilistic RNA-seq
538 quantification. *Nature Biotechnology*, 34(5):525–527, 2016.
- 539 [71] Heng Li and Richard Durbin. Fast and accurate long-read alignment with Burrows–Wheeler trans-
540 form. *Bioinformatics*, 26(5):589–595, 2010.
- 541 [72] Doug Hyatt, Gwo-Liang Chen, Philip F. LoCascio, Miriam L. Land, Frank W. Larimer, and Loren J.
542 Hauser. Prodigal: prokaryotic gene recognition and translation initiation site identification. *BMC
543 Bioinformatics*, 11(1):119, March 2010.
- 544 [73] Martin Steinegger and Johannes Söding. Clustering huge protein sequence sets in linear time. *Nature
545 Communications*, 9(1):2542, 2018.
- 546 [74] Maria Hauser, Martin Steinegger, and Johannes Söding. MMseqs software suite for fast and deep
547 clustering and searching of large protein sequence sets. *Bioinformatics*, 32(9):1323–1330, 01 2016.
- 548 [75] Rachel Seongeon Kim, Eli Levy Karin, Milot Mirdita, Rayan Chikhi, and Martin Steinegger. Bfd—a
549 large repository of predicted viral protein structures. *Nucleic Acids Research*, 53(D1):D340–D347,
550 11 2024.
- 551 [76] Fabian Sievers, Andreas Wilm, David Dineen, Toby J Gibson, Kevin Karplus, Weizhong Li, Rodrigo
552 Lopez, Hamish McWilliam, Michael Remmert, Johannes Söding, Julie D Thompson, and Desmond G
553 Higgins. Fast, scalable generation of high-quality protein multiple sequence alignments using clustal
554 omega. *Molecular Systems Biology*, 7(1), January 2011.
- 555 [77] Martin Steinegger and Johannes Söding. MMseqs2 enables sensitive protein sequence searching for
556 the analysis of massive data sets. *Nature Biotechnology*, 35(11):1026–1028, November 2017. Publisher:
557 Nature Publishing Group.
- 558 [78] Antonio Pedro Camargo, Simon Roux, Frederik Schulz, Michal Babinski, Yan Xu, Bin Hu, Patrick
559 S. G. Chain, Stephen Nayfach, and Nikos C. Kyrpides. Identification of mobile genetic elements with
560 geNomad. *Nature Biotechnology*, 42(8):1303–1312, August 2024.
- 561 [79] Mosè Manni, Matthew R Berkeley, Mathieu Seppey, Felipe A Simão, and Evgeny M Zdobnov.
562 BUSCO Update: Novel and Streamlined Workflows along with Broader and Deeper Phylogenetic
563 Coverage for Scoring of Eukaryotic, Prokaryotic, and Viral Genomes. *Molecular Biology and Evolu-*
564 *tion*, 38(10):4647–4654, September 2021.
- 565 [80] Evgenia V Kriventseva, Dmitry Kuznetsov, Fredrik Tegenfeldt, Mosè Manni, Renata Dias, Felipe A
566 Simão, and Evgeny M Zdobnov. OrthoDB v10: sampling the diversity of animal, plant, fungal,
567 protist, bacterial and viral genomes for evolutionary and functional annotations of orthologs. *Nucleic
568 Acids Research*, 47(D1):D807–D811, January 2019.
- 569 [81] Martin Larralde. Pyrodigal: Python bindings and interface to Prodigal, an efficient method for gene
570 prediction in prokaryotes. *Journal of Open Source Software*, 7(72):4296, April 2022.

- 571 [82] Martin Larralde and Georg Zeller. PyHMMER: a Python library binding to HMMER for efficient
572 sequence analysis. *Bioinformatics*, 39(5):btad214, May 2023.
- 573 [83] James Robertson and John H. E. Nash. MOB-suite: software tools for clustering, reconstruction and
574 typing of plasmids from draft assemblies. *Microbial Genomics*, 4(8), August 2018.
- 575 [84] Michael Feldgarden, Vyacheslav Brover, Narjol Gonzalez-Escalona, Jonathan G Frye, Julie Haendiges,
576 Daniel H Haft, Maria Hoffmann, James B Pettengill, Arjun B Prasad, Glenn E Tillman, et al.
577 AMRFinderPlus and the Reference Gene Catalog facilitate examination of the genomic links among
578 antimicrobial resistance, stress response, and virulence. *Scientific reports*, 11(1):12728, 2021.
- 579 [85] Célio Dias Santos-Junior, Shaojun Pan, Xing-Ming Zhao, and Luis Pedro Coelho. Macrel: antimicrobrial
580 peptide screening in genomes and metagenomes. *PeerJ*, 8:e10555, 2020.
- 581 [86] Stephen Nayfach, Antonio Pedro Camargo, Frederik Schulz, Emiley Eloë-Fadrosh, Simon Roux, and
582 Nikos C. Kyrpides. CheckV assesses the quality and completeness of metagenome-assembled viral
583 genomes. *Nature Biotechnology*, 39(5):578–585, May 2021.
- 584 [87] V. A. Traag, L. Waltman, and N. J. Van Eck. From Louvain to Leiden: guaranteeing well-connected
585 communities. *Scientific Reports*, 9(1):5233, March 2019.
- 586 [88] Salvador Capella-Gutiérrez, José M Silla-Martínez, and Toni Gabaldón. trimAl: a tool for automated
587 alignment trimming in large-scale phylogenetic analyses. *Bioinformatics*, 25(15):1972–1973, 2009.
- 588 [89] Morgan N. Price, Paramvir S. Dehal, and Adam P. Arkin. FastTree 2 – Approximately Maximum-
589 Likelihood Trees for Large Alignments. *PLoS ONE*, 5(3):e9490, March 2010.
- 590 [90] Luiz Irber, N Tessa Pierce-Ward, Mohamed Abuelanin, Harriet Alexander, Abhishek Anant, Keya
591 Barve, Colton Baumler, Olga Botvinnik, Phillip Brooks, Daniel Dsouza, et al. sourmash v4: A
592 multitool to quickly search, compare, and analyze genomic and metagenomic data sets. *Journal of
593 Open Source Software*, 9(98):6830, 2024.
- 594 [91] Mark Raasveldt and Hannes Mühleisen. Duckdb: an embeddable analytical database. In *Proceedings
595 of the 2019 international conference on management of data*, pages 1981–1984, 2019.

596 3 Methods

597 3.1 Performance-optimized and cloud-accelerated genome assembly of all SRA ac- 598 cessions

599 We designed a cloud architecture to perform SRA-wide genome assembly efficiently and in parallel. Fig-
600 ure 1b describes the workflow. Briefly, a Docker container processes each SRA accession independently:
601 1) raw reads of an accession are downloaded from a cloud mirror of the SRA, 2) a conservative assembly
602 (unitigs) of the accession is made from the reads, 3) a consensus assembly (contigs) is made from the
603 unitigs, and 4) both the unitigs and contigs are compressed and uploaded to a public repository.

604 We executed this system at SRA-scale using cloud resources. Containers are executed in parallel over
605 tens of thousands of cloud computers through a container orchestration system, and a set of dashboards
606 were deployed to monitor the execution. Extended Data Fig. 1 shows some key statistics of the execution.
607 We optimized the system to use as many CPU cores in parallel as possible, as opposed to running it
608 at smaller scale over a longer period of time, to take advantage of lower computation costs and higher
609 availability of cloud instances during night time.

610 Using this system we performed genome assembly over the entire SRA and report results for each
611 accession in two forms: unitigs (non-branching paths of the reads de Bruijn graph, here $k = 31$) and contigs
612 (non-branching paths of the de Bruijn graph after simplification steps). The rationale for this dual set of
613 results is to allows researchers to choose between a near-lossless but fragmented representation (unitigs)
614 and a more lossy but more contiguous one (contigs). Unitigs are shorter sequences that preserve almost
615 all the genetic variation from the original reads, including alternate alleles involving single nucleotide or
616 indel variants. Unitigs are ideal for sensitive, k -mer-based searches where finding even minor variants is
617 critical. Contigs, on the other hand, are the more contiguous, consensus assemblies. To obtain contigs,
618 the assembler extends the unitigs by collapsing biological variations and removing any remaining putative
619 technical sequencing error. It makes a “best guess” of the errors to remove and major allele to collapse into
620 a single, representative path. Contigs are optimized for downstream applications like protein prediction
621 and gene identification, where sequence contiguity is more important than preserving minor variants.

Dataset	Type	Reads	# accessions	Index/Seqs	Compress.
Sourmash Branchwater	Index	1.4 Pbp	682,688	7 TB	191×
MetaGraph-SRA	Index	3.3 Pbp	1,891,328	7 TB	473×
Pebblescout-SRA	Index	3.7 Pbp	4,141,058	171 TB	21×
NCBI SRA (Dec 2023)	Raw data	50.3 Pbp	27,764,168	19 PB	3×
Logan, Unitigs (v1)	Assembly	48.3 Pbp	27,311,279	2 PB	24×
Logan, Contigs (v1.1)	Assembly	44.1 Pbp	26,788,829	315 TB	140×
Logan-Search	Index	43.9 Pbp	23,404,655	1 PB	44×

Table 1: **Size of existing indexed data vs Logan.** The MetaGraph-SRA and Pebblescout-SRA rows refer to all SRA accessions indexed by MetaGraph and Pebblescout respectively [4, 6]. The NCBI SRA row refers to all public accessions from the SRA as of December 10th 2023. The Reads column refers to the number of bases in SRA reads for the considered dataset. The Index/Seqs column indicates the sum of all sub-indices sizes (for Branchwater [8], MetaGraph and Pebblescout) or the size of compressed sequences for all accessions (for SRA and Logan). The “Compress.” column gives the compression ratio between the size of reads (Reads column) as if each base was stored using 8 bits, and the Index/Seqs column.

622 3.1.1 Input data

623 We selected all public samples from the Sequence Read Archive on December 10th, 2023, with read length
624 above 31 bp. Accessions with shorter reads than 31 bp would yield no usable k -mers in the downstream
625 assembly step. The list of accessions was obtained from the NCBI SRA metadata table, using Amazon Web
626 Services (AWS) Athena with SQL query filter: WHERE consent = 'public' and avgspotlen >= 31.
627 This resulted in 27,764,168 accessions totalling 50,304,659,857 bases in reads. In the rest of this manuscript
628 we refer to this dataset as 'the SRA', although the current-day SRA has since been updated with new
629 samples.

630 3.1.2 Assembly tools

631 Unitigs were constructed using a modified version of Cuttlefish2 [38] (commit 9401ef5 of forked repository
632 github.com/rchikhi/cuttlefish), augmented to record approximate mean k -mer abundance per
633 unitig. We also modified the k -mer counting method KMC3 [39] integrated inside of Cuttlefish2 to
634 stream SRA files directly through a piped call to fasterq-dump with parameters --seq-defline '>'
635 --fasta-unsorted --stdout, avoiding a prior decompression step to disk, and discarding on the fly
636 FASTQ headers and quality values.

637 The original version of Cuttlefish2 did not record any abundance value. We modified Cuttlefish2 to
638 record abundances values per k -mer during construction, then to report the average abundance over all
639 k -mers of a unitig. However, the reported per- k -mer abundances were approximated with two heuristics:
640 1) due to a technicality during graph construction, abundances of $k + 1$ -mers were recorded, hence the
641 abundance of each k -mer was obtained by summing the abundances of all the $k + 1$ -mers it appeared in,
642 then dividing the sum by two. 2) To save memory during graph construction, abundances were stored
643 in a 8 bits encoding scheme that maintains an error of not more than 5%, with a maximum abundance
644 value of 50,000. To remove some of the likely sequencing errors, k -mers seen only once in an accession
645 were discarded from unitigs.

646 Unitigs were given as input to Minia3 [40] (commit 71484e8 of github.com/GATB/minia), which per-
647 forms de Bruijn graph simplifications and error-correction following closely the heuristics designed by
648 the SPAdes assembler [41], a tool well-established for its high accuracy and low rate of misassemblies in
649 metagenomic contexts. We expect Minia3 contigs to share a similar accuracy profile, since the modifi-
650 cations focused primarily on performance optimization and memory frugality, rather than altering the
651 fundamental algorithms responsible for assembly quality. Minia3's accuracy was independently validated
652 in [42, 43, 44]. The resulting contigs that are longer than 150 bp and are connected to at least one other
653 contig were reported. For both unitigs and contigs, the FASTA headers contain link information (in the
654 format of BCALM2 [45] output) that enable to reconstruct the assembly graph in GFA format. In the
655 unitigs and contigs file of an accessions all 31-mers are distinct, by construction.

656 These tools were selected for their memory and running time frugality. In addition, Minia3 was
657 chosen for also its conservative approach to graph simplification and hence higher retention of sequence
658 complexity. A comparison with two other state-of-the-art assemblers (Penguin [46], rnaviralSPAdes [47])
659 is provided in Extended Data Fig. 1, showing that substantial cloud computing costs were saved by this
660 pipeline.

661 To compress unitigs and contigs, we developed and applied a novel block variant of the Zstandard
662 algorithm [48] (<https://github.com/as1/f2sz>). The compressor creates FASTA-aligned blocks allowing
663 for faster random access to any subset of contigs. It remains compatible with the ubiquitous `zstd -d`
664 and `zstdcat` decompression command line tools.

665 3.1.3 Cloud infrastructure

666 For SRA-scale assembly we have set up a cloud infrastructure on Amazon Web Services (AWS), to
667 perform assembly of each SRA accession in a fully parallel fashion. In brief, the infrastructure is based

668 on a Docker container executing a set of Python scripts responsible for calling child programs for unitig
669 and contig assembly and validation. The AWS Batch execution system handles scheduling of containers
670 across a pool of cloud computers (AWS EC2 instances). Each container is executed independently for
671 each SRA accession. Each instance is equipped with temporary network storage (EBS). The number of
672 CPUs, RAM, and storage for each job were set according to size of input reads measured in megabases.
673 The Batch instances pool was set to the `c6g` and `c7g` families (AWS Graviton-based instances), with sizes
674 `4xlarge` and above to target larger instances and thus limit the number of instance creation API calls.

675 Executions were monitored primarily through live dashboards on AWS CloudWatch and Batch ser-
676 vices, as well as a DynamoDB database recording runtime and assembly statistics for each accession.
677 Global statistics such as number of processed accessions and total size of raw data assembled were recorded
678 in real-time by sending messages from each container to a global partitioned database. Summary metrics
679 were aggregated, enabling to monitor computation speed during execution and potentially stop all jobs
680 should metrics fall behind projected estimates.

681 To limit the number of simultaneous queries to NCBI servers, two mechanisms were implemented:
682 (1) raw `.sra` files were directly downloaded from the AWS Registry of Open Data cloud mirror of the
683 NCBI SRA to a cloud instance in the same data center (`us-east-1`), and (2) special `.sra` files containing
684 alignments to RefSeq were handled by downloading references from a S3 mirror of RefSeq. Aligned `.sra`
685 file containing references from the NCBI WGS database were discarded in the later runs, but some were
686 processed in the earlier runs before we identified that they incurred (rate-limited) queries to NCBI servers.

687 Results (unitigs, contigs) have been deposited in a public repository. Detailed instructions to download
688 the data are provided here: <https://github.com/IndexThePlanet/Logan>.

689 While the cloud infrastructure used to construct Logan is solely intended for internal usage due to
690 its high execution costs, its source code is publicly available at <https://gitlab.pasteur.fr/rchikhi-pasteur/erc-unitigs-prod/>.

692 3.1.4 Assembly results

693 In total 27.3 million accessions were assembled into unitigs, representing 96% of the SRA in size as of
694 December 2023. Some accessions resulted in too many unitigs to fit the assembly graph in memory, hence
695 were not further assembled into contigs at this time. 26.8 million accessions were assembled from unitigs
696 to contigs, representing 88% of the SRA in size. The total cloud computation time (unitigs and contigs)
697 was around 30 million CPU hours (1).

698 Assembly contiguity statistics for the Logan contigs are reported in Extended Data Fig. 1. Contigs
699 for Whole-Genome Sequencing/Amplification (WGS/WGA) accessions are generally longer than those of
700 RNA-Seq accessions or other sequencing types, as expected by the longer sequenced molecules. Note that
701 non-circular contigs shorter than 150 bp that are isolated nodes in the assembly graph were discarded by
702 the assembler as they were more likely to be artifacts than actual biological material.

703 Across Logan, standard assembly metrics were computed using seqkit [49] and stored into a database
704 (number of unitigs, contigs, N50 values, total length of assemblies, longest assembled sequence per acces-
705 sion). In addition, FASTA file sizes before and after Zstandard compression were also recorded. All these
706 statistics were stored on a AWS DynamoDB database then exported to a public repository (S3 bucket) in
707 the Parquet format (<https://github.com/IndexThePlanet/Logan/blob/main/Stats-v1.md>), enabling
708 users to link this database with other NCBI SRA databases such as STAT [5] or SRA metadata.

709 3.2 Shallow and deep homology search in Logan contigs

710 For translated-protein searches over all Logan contigs, DIAMOND2 v2.1.9 [11] was run with parameters
711 `-b 0.4 --masking 0 -s 1 --sensitive` to balance speed and sensitivity. The query sequence(s) were
712 provided as indexed references to DIAMOND2, and Logan contigs were streamed accession by accession
713 as queries. For nucleotide searches over all Logan contigs, minimap2 v2.28 [50] was run with parame-
714 ters `-x sr --sam-hit-only -a` to return all contig sequences and their alignment, optimizing for short

715 matches. A custom cloud pipeline using AWS Batch was set up to perform those searches, which each
716 take approximately 11 hours on 60,000 vCPUs. The pipeline and infrastructure source code is available
717 at https://gitlab.pasteur.fr/rchikhi_pasteur/logan-analysis.

718 3.3 SRA-wide public search engine

719 We developed Logan-Search, a publicly available search engine able to approximately locate a queried
720 DNA sequence among 23.4 million accessions. Logan-Search performs k -mer-based queries. A k -mer is a
721 word of length k ($k = 31$ is used in this context). Within minutes, Logan-Search identifies the accessions
722 to which each k -mer from the queried sequence is associated. This enables us to provide a similarity
723 metric between the query and each indexed accession.

724 Constructing the Logan-Search engine required to index all k -mers from Logan unitigs, in order to offer
725 a way to instantaneously detect to which sample a k -mer belongs to. We built this index using kmtricks [51]
726 (kmtricks-logan tag at github.com/tlemane/kmtricks) and kmindex [52] (v0.5.3). Computations
727 were performed on Microsoft Azure cloud platform, using an Azure Batch Pool consisting of 625 virtual
728 machines (VMs) of type Standard_D32d_v5. The workload was divided into approximately 45,000 tasks
729 used to construct partial indexes, which were subsequently merged to produce the final index. A small
730 fraction of tasks requiring larger memory capacity were executed on Standard_D96d_v5 instances. In
731 total, the computation took around 10 days to complete.

732 In total, approximately 2×10^{15} k -mers were indexed. The accessions were separated into groups
733 on the basis of their library source (e.g. genomic, transcriptomic, metagenomic, metatranscriptomic,
734 etc...) and of their superkingdom phylogeny classification, obtained from the NCBI STAT database [5].
735 Exceptions were made for human and mouse accessions, which were classified separately. The kmindex
736 tool builds Bloom Filters [53], with a non-null false positive rate, tuned to be approximately 0.005% by
737 setting parameters adapted to the size of each indexed dataset and by using the Findere algorithm [54].
738 The overall size of the index is approximately one petabyte, stored on disk. At query time, only specific
739 target sub-parts of the index are mapped on RAM. For a thousand basepair sequence query, results are
740 obtained in approximately 6 minutes, using 12 small 4-vCPUs virtual machines.

741 The kmindex tool retrieves accessions containing sequences similar to a query, based on the percentage
742 of k -mers from the query existing in the accession. On top of those results, we developed a visualization
743 interface based on kmviz (<https://github.com/tlemane/kmviz>) that offers several features:

- 744 • All metadata from SRA associated to each accession is made available. The interface enables
745 researchers to visually retrieve those metadata: geographic localization of accessions on a worldwide
746 map, and drawing highly tunable plots from discrete or textual attributes.
- 747 • A summary of the query is generated using a large language model (GPT-4o) [55], based on the
748 SRA metadata associated to top hits, allowing users to quickly assess potentially relevant contextual
749 information, such as organism or location.
- 750 • Logan-Search natively returns a list of accessions that contain the query sequence, but it does not
751 identify the specific sequence within each accession that matches the query. To fill this gap, a
752 microservice based on the back_to_sequences tool [56] is used to retrieve, on a per-accession basis,
753 the contig or unitig sequences that match the query. Additionally, a BLASTn [57] alignment using
754 default parameters is performed between the query and the extracted contigs or unitigs; it completes
755 instantly as both aligned sequences are typically kilobase-sized.

756 3.4 Pilot plastic active enzyme search

757 In a pilot experiment to assess the feasibility of a plastic-active enzyme search, we retrieved 102 publicly
758 available, experimentally validated, polyethylene terephthalate (PET)-degrading A/B hydrolases from
759 the PAZy database (sequence accessed Aug 26, 2024 from GenBank or PDB accessions [18]). These

reference sequences were supplemented with 12 previously computationally identified PETases and 74 MGnify sequences with predicted structures similar to *IsPETase*, found with Foldseek [58, 59, 13]. To reduce redundancy, these 188 enzymes were clustered at 90% aaid ('USEARCH v11.0.667_i86linux32', **-cluster_fast -id 0.90**), resulting in 153 representative sequences [60].

We focused on 56 sequences related to *IsPETase* by at least 30% aaid (USEARCH v11.0.667_i86linux32, **-cluster_fast -id 0.30**) [60]. We filtered Logan retrieved contigs to " *IsPETase-like*" sequences with high identity and confidence, resulting in 230,804 hits ('DIAMOND2 blastx v2.1.9' as above; e-value < 1e⁻⁸ and aaid > 40%) [11].

As we were interested in potentially active enzymes with intact catalytic cores, we generated a multiple structure alignment of *IsPETase*-like reference sequences. Structures were predicted using AlphaFold3 and aligned with Muscle3D (v5.1.linux64, **-align**) [61, 62]. The amino acids of the subalignment containing the conserved core were extracted manually and were used to generate a custom HMM ('HMMER v3.4'; **esl-reformat stockholm, hmmbuild**) [63]. Next, we identified stop-stop ORFs ('EMBOSS v6.6.0.0'; **getorf**) and filtered for those with at least 90% coverage of the core HMM ('HMMER v3.4'; **hmmscan**), resulting in 2,272 unique sequences with a maximum e-value of 3.8e⁻⁴³.

To infer the evolutionary relationships of these Logan hits, we clustered them at 95% aaid (853 representative sequences), and used IQ-TREE2 with 1000 bootstrap iterations to generate a tree (v2.4.0, **-b 1000**) [64]. Ancestral reconstruction of these sequences was performed with GRASP command-line [65], resulting in 175 ancestral sequences. For initial synthesis and activity testing, 161 candidate sequences from Logan were manually chosen over the tree to encompass an even sequence diversity. Twenty-one ancestrally reconstructed sequences were included, as well as six distantly related sequences from MGnify as expected negative controls, and 11 known PETases as positive controls, amounting to 199 sequences in total.

After the PETase activity was measured for the initial round of synthesized sequences (see Section 3.8) the phylogenetic tree was re-sampled in areas of relatively enriched PETase activity through manual identification of closely related and sister sequences (n = 13) for additional *in-vitro* screening.

3.5 Plastic active enzyme homolog expansion

As the initial search seed, we used the PAZy database of experimentally validated plastic-active enzymes (Accessed Dec 19, 2024) [18]. All publicly available PAZy sequences were retrieved from GenBank by accession (213/245, 86.94%). To generate a PAZy phylogenetic network (Extended Data Fig. 2a), we performed all-vs-all alignment (usearch v11.0.667) between sequences (213 nodes), and retained significant alignments (3,195 edges, ≥30% amino acid identity and e-value < 1e⁻⁵). The sequences clustered into 42 graph components, of 11 structurally distinct protein folds (CATHdb [21]), with minimum-spanning tree edges shown. The Alpha/Beta Hydrolase protein family are the most represented, with 170/213 (79.8%) of the sequences and 28/42 (66.7%) of the components.

Notably, Alpha/Beta hydrolases, amidases, beta-lactamase, and arylesterase protein-fold families have reported activity for more than one plastic substrate, and 16 individual sequences have reported activity for more than one plastic substrate. Thus plastic degrading activity is a polyphyletic trait, and suggests it may be an incidental, or off-target function of these enzymes, which broadly function to hydrolyze organic molecules such as lipids, esters, or carbohydrates. This suggests these enzymes and their homologs may contain additional plastic substrate activity.

To identify PAZy homologs, we queried the sequences into the NCBI non-redundant protein database ('nr', accessed: 2024-12-27) [10] using DIAMOND2 (**--very-sensitive**) [11] which retrieved 5,593,290 unique sequences (e-value < 1e⁻⁵). Of these, 2,735,790 sequences contained a HMM match ('hmmscan'; < 1e⁻⁵ and 95% model coverage) against a PAZy protein domain (Pfam models: PF00082, PF00089, PF00144, PF00561, PF01083, PF01425, PF01522, PF01674, PF01738, PF02983, PF03403, PF03576, PF06850, PF07224, PF07732, PF09995, PF10503, PF10605, PF12146, PF12695, PF13472, PF20434, PF21419, PF24708). Incidentally, we did not include a model for Cyclase-like domains, which resulted in

808 zero novel sequences being called.

809 Next, we queried these centroids against our newly created Logan assemblage (26.8 million datasets,
810 50.0 Pbp, 11.5 hours), as per Section 3.2, and retrieved an additional 6.5 billion hits, (DIAMOND, evalue
811 < $1e^{-8}$), which were filtered as above.

812 3.6 Yeast strains and PETase expression vectors

813 All yeast strains generated in this study are derivatives of DHY213 [*MATa CAT5(91M) SAL1 MIP1(661T)*
814 *HAP1 MKT1(30G) RME1(INS-308A) TAO3(1493Q) leu2Δ0 his3Δ1 ura3Δ0 met15Δ0*], a modified ver-
815 sion of BY4741 [66] with increased sporulation efficiency, mitochondrial stability and efficient biosynthetic
816 gene expression [67].

817 PETase encoding genes were expressed as surface displayed or secreted constructs from episomal
818 plasmids derived from the pJC170 backbone [68] (gift from Jef Boeke) modified to contain an additional
819 marker (*natMX*) and the surface display or secretion expression cassettes. For surface display, PETase
820 encoding genes were fused with the Ccw12 secretion signal (amino acid 1 to 19) on their 5' ends and to
821 a flexible [AGSAGSAAGSG] linker, a Myc tag and the remainder of the Ccw12 amino sequence (amino
822 22 to 133) on their 3' ends. For secretion, the architecture of the construct was the same as for surface
823 display but without the cell wall anchoring domain of Ccw12 (amino 22 to 133). Finally, the empty
824 vectors additionally contain a placeholder sequence with two inverted BsaI restriction sites between the
825 Ccw12 secretion signal and the linker sequence in order to facilitate backbone cleavage for recombination
826 or golden gate cloning. The sequence for these plasmids (pRLK152 for surface display; pRLK153 for
827 secretion) is provided in the supplemental material.

828 3.7 Synthetic DNA cloning

829 PETase encoding genes were synthesized by Twist Biosciences (USA) and contained the standard Twist
830 adapters and 40 base pairs of homology on the 5' (5' CGCTTCTATCGCCGCTGTCGCAGCTGTCGCTTCT-
831 GCCGCA) and 3' (5' GCGGGTTCTGCTGCTTCTGCTGCTGGTTCTGGTGAATTG) ends to fa-
832 cilitate *in vivo* recombination in yeast. Between 10 and 20 ng of synthetic DNA was transformed in
833 yeast using the standard lithium acetate method along with 25-50 ng of plasmid backbone (pRLK152
834 or pRLK153). Yeast transformants were selected on synthetic medium containing clonNAT antibiotic
835 and lacking uracil (SD/MSG-ura+ clonNAT: 1.7 g/L yeast nitrogen base without amino acids without
836 ammonium sulfate, 1 g/L monosodium glutamate, 20 g/L dextrose, 20 g/L agar, 100 µg/mL clonNAT).
837 The pool of transformants for each construct was maintained as a single colony on a colony array.

838 3.8 BHET halo assay

839 To prepare assay plates, a 1 M BHET (CAS# 959-26-2, Sigma-Aldrich) solution was prepared by diluting
840 BHET flakes into 100% DMSO and heating slightly until complete dissolution. The BHET solution
841 was then added to YPD medium (10 g/L yeast extract, 20 g/L peptone, 20 g/L dextrose, 20 g/L agar,
842 100 µg/mL clonNAT) prior to pouring Omnitray plates (ThermoFisher). Plates were kept at room
843 temperature to prevent BHET recrystallization.

844 Yeast strains expressing surface displayed or secreted *IsPETase* were pinned in 384-array format (4
845 colonies per construct) onto YPD+BHET plates using a colony pinning robot (Singer Instruments, United
846 Kingdom) and incubated at 30°C for up to 5 days. After incubation, yeast colonies were washed off the
847 plates using water and a cell spreader tool to reveal a white halo or a clearing (loss of opacity). Each plate
848 was imaged before robotic pinning as well as before and after colony washing (spImager, SP Robotics Inc,
849 Canada).

850 BHET halo measurement was implemented in R as follows: First colonies were identified using the
851 gitter package and the resulting colony mask was applied to the plate image before pinning (background
852 image) and after colony washing (halo image) to extract pixel coordinates corresponding to each colony

area in the image. Pixel intensity was extracted from each image using the imager package and the median pixel intensity for each colony area was then determined in the halo image and the median pixel intensity for the same colony area in the background image was subtracted. To normalize pixel intensity values across plates, the average median pixel intensity obtained across 4 replicate colonies containing the empty surface display or secretion plasmids present on each plate was subtracted from all colony median pixel intensity values on the given plate.

3.9 High performance liquid chromatography

Isogenic clones were isolated from each colony with activity in the BHET halo plate assay and grown to saturation in YPD containing 100 µg/mL cloNAT at 30°C. The saturated culture was then diluted 1000-fold in YPD+cloNAT and grown for 24 hours at 30°C. 95 µl of culture was transferred into a 96-well plate prior to adding 5 µl of 500 mM or 250 mM BHET in 100 % DMSO. After 17 hours of incubation at 30°C, an aliquot of each reaction was diluted 10 or 20 times in 100% DMSO for reactions in 12.5 and 25 mM BHET respectively, centrifugated for 2 minutes at 3000 rpm and the supernatant was stored at -20°C.

Supernatants were fractionated on reversed-phase HPLC using an HP1050 system (HP/Agilent, USA) mounted with a Zorbax SB-C8 column (4.6 x 150 mm, 5 µm). The column was maintained at 22-24°C. Analytes were eluted over 37 minutes with 0.1% formic acid in water (aqueous solvent) and 0.1% formic acid in acetonitrile (organic solvent) using the following gradients: 1 to 5% organic (vol/vol) over 20 minutes at 0.8 ml/min, 5 to 52.5% organic (vol/vol) over 14 minutes at 0.8 ml/min, 52.5% to 100% organic (vol/vol) and 0.8 to 3.0 ml/min over 0.2 min, 100% organic (vol/vol) for 0.8 min at 3.0 ml/min, 100% to 1% organic (vol/vol) and 3.0 to 0.8 ml/min over 1.0 min, and 1% organic (vol/vol) for 0.2 minutes. Detection wavelength was 240 nm with a 4 nm bandwidth. Peak identities were established using commercial TPA ≥98% purity (CAS: 100-21-0, Sigma-Aldrich), MHET ≥95% purity (CAS: 1137-99-1, Advanced ChemBlocks), and BHET ≥95% purity (CAS: 959-26-2; Sigma-Aldrich). Data analysis was performed in R using the chromatographR package and analyte abundance was determined by measuring absorbance peak area at 240 nm. TPA and its ester derivatives (BHET and MHET) have similar extinction coefficients at 240-244 nm [69] and purified BHET dimer (see below) consistently gave peak areas that were 3.5 times smaller than pure BHET across a range of analyzed amounts (Extended Data Fig. 3 f-g). Therefore, product formation was expressed as a ratio between the peak area of the enzymatic reaction products TPA, MHET and 3.5-times the peak area of BHET dimer, relative to the sum of all peaks (TPA, MHET, 3.5 times BHET dimer, and BHET). Finally, the relative abundance of each analyte was normalized to the cell concentration at the time of BHET addition in each reaction determined using a Beckman-Coulter Counter Z1 equipped with a 100 micron aperture tube.

For BHET conversion measurement from halo zones on YPD+BHET agar plates, colonies were washed off of the plates and white halos or clearings were excised using a pipet tip with a 1-2 mm diameter. The agar plugs were soaked in 250 ul of 100% DMSO for 24 hours at room temperature prior to centrifugation at 15,000 rpm for 3 minutes. The supernatant fraction was further diluted 5-fold prior to HPLC analysis (see above) or LC-MS (see below). Samples for LC-MS analysis were prepared in low-bind tubes and centrifuged for 20 minutes at 15,000 rpm prior to analysis to prevent any transfer of large particles into the instrument.

3.10 BHET dimer purification

20 µmoles of BHET from a 1M BHET stock solution (CAS: 959-26-2, Sigma-Aldrich) which contained BHET dimer (O,O' -(ethane-1,2-diyl) bis(oxy(2-hydroxyethyl)carbonyl)terephthalate) as impurity was fractionated by HPLC using the aforementioned fractionation protocol. Fractions between 28-29 minutes and 33.9-35 minutes were collected from the waste line of the instrument prior to being dried at 60°C for 24 hours to isolate BHET and BHET dimer, respectively. The weight of the dried fraction was then measured on a high-precision scale and the dried fraction was resuspend in 100% DMSO to a final

900 concentration of 1 M. Fraction purity was verified by HPLC analysis across a range of concentrations (1
901 mM – 125 μ M, Extended Data Fig. 3e-f).

902 3.11 Liquid chromatography - Mass spectrometry analysis

903 Halo zones from under yeast colonies were extracted from agar plates as described above. An Agilent
904 1260 Infinity II with 6545 LC/QTOF mass spectrometer was used to analyze the samples in positive
905 ionization mode with Dual AJS electrospray ionization (ESI) equipped with Agilent ZORBAX Eclipse
906 Plus C18 column (2.1x50mm, 1.8- μ m particles) and ZORBAX Eclipse Plus C18 guard column (2.1x5mm,
907 1.8- μ m particles). LC parameters were as follows: injection volume 2 μ L preceded with a 4 μ L needle wash
908 with sample, autosampler chamber temperature 20°C, column oven temperature 40°C. Mass spectrometry
909 parameters were as follows: gas temperature 320°C, drying gas flow 10 L/min, nebulizer 35 psi, sheath
910 gas 350°C at 11 liters per minute, VCap 3500V, Nozzle voltage 1000V, fragmentor 125V, skimmer 65V.
911 The solvent gradient with a flow of 0.5 ml per minute started with 99% mobile phase A (Optima LC/MS
912 H₂O+0.1% Formic Acid, Fisher Chemical P/N LS118-4) and 1% mobile phase B (Optima LC/MS Ace-
913 tonitrile+0.1% Formic Acid, Fisher Chemical P/N LS120-4), kept for five minutes, increased linearly to
914 100% B at 5 minutes, followed by five minutes at 100% B, then back to 1% B over 2 minutes and finally
915 held at 1% B for an additional 5 minutes. The post-run time was two minutes (instrument conditioning
916 at 99% mobile phase A). The raw data was analyzed using Agilent MassHunter Qualitative Analysis 12.0.
917 Counts of molecules with mass-to-charge (m/z) ratios specific to MHET, BHET and BHET dimer were
918 collected, and the area under the curve of each peak was calculated to determine the abundance of each
919 molecule. MHET and BHET dimer abundance was expressed relative to spectral counts obtained for
920 BHET in each sample. A putative BHET dimer structure was inferred from the m/z ratio.

921 3.12 Discovery and characterization of HHV-6 reactivation

922 For all HHV-6 analyses, we utilized Logan-Search using the AF157706 reference genome and transcriptome,
923 reflecting the HHV-6B strain, which is endemic outside of Sub-Saharan Africa. For selecting HHV-6
924 transcripts to query in Logan-Search, we prioritized the HHV-6 genes annotated as ‘late’ that encode proteins
925 essential for viral assembly and release of particles. In doing so, our search prioritized libraries
926 containing viral expression consistent with full reactivation (rather than latency or early reactivation).
927 The full RefSeq transcripts for U83 and U91 were input to Logan-Search and queried against the full set
928 of human RNA-Seq datasets. Once these accessions were identified from the full search, individual SRR
929 files were downloaded and further analyzed for modality-specific analyses. For bulk RNA-seq libraries,
930 **kallisto**[70] in quant mode was used with the HHV-6 transcriptome as previously described [22]. For
931 single-cell analyses of the organoid system, processed human counts matrices were downloaded from GEO
932 and further annotated with the **kallisto bus** single-cell counts for HHV-6. For ChIP-seq libraries, raw
933 .fastq files were re-mapped to the HHV-6 reference genome using bwa [71] with downstream analyses
934 conducted using **GenomicAlignments**. Metadata, including the clinical trial identifier, input cell type,
935 and donor identity, was pulled from the SRA metadata annotations per BioProject.

936 To mitigate sources of confounding for downstream analyses (i.e., errant, non-HHV-6 detection), we
937 performed a series of stringent filters for read quantification. First, to minimize the possibility of multi-
938 mapping transcripts, we excluded the HHV-6 DR1 transcript that possesses high homology with human
939 transcripts [22]. Second, to mitigate the possibility of HHV-6 contamination (either environmental or
940 index hopping), we a limit of detection per library requiring (a) a minimum of 3 unique genes and (b)
941 minimum of 10 unique sequencing reads to call a library positive. Finally, for libraries with high viral
942 reactivation, we verified single nucleotide diversity comparing mutations with a minimum cover of 10x
943 per library and allele frequencies exceeding 90% in one library and less than 10% in the other Extended
944 Data Fig. 6d). Libraries meeting these criteria were further processed to estimate the viral RNA counts
945 per million defined as the non-DR1 HHV-6 transcripts divided by the total sequencing reads per library.

946 Hence, we emphasize that these additional measures yield a conservative estimate of HHV-6 abundance
947 in these libraries.

948 3.13 Gene calling and protein clustering

949 Protein-coding genes were predicted in all assembled Logan contigs using Prodigal [72] (version 2.6.3) in
950 metagenomic mode (`-p meta`). The predicted proteins were then divided into 'human' vs. 'other' based
951 on SRA metadata associated with their contigs through their SRA accessions and into 'complete' vs.
952 'partial', based on Prodigal's output. In all, three subsets were obtained: 31.2 billion "human-complete",
953 109.4 billion "other-complete", and 304.7 billion "human-partial" protein sequences. We excluded the
954 "other-partial" set because partial predictions are more error-prone and because this set alone contained
955 nearly one trillion sequences that did not cluster well.

956 These subsets were then clustered using Linclust [73] (commit `62a2ad`) on AWS Batch/EC2; the
957 input was streamed from Amazon S3 and partitioned into fixed-size line chunks sized to instance memory
958 (up to 1.5 billion proteins per chunk). Jobs ran as Batch array jobs on x2gd.metal nodes with a 700G
959 split-memory limit (`--split-memory-limit 700G`).

960 Linclust was run in several cascaded rounds [74]. First, we removed near-duplicate protein fragments
961 at 90% sequence identity (`--min-seq-id 0.9`) with target coverage 90% (`--cov-mode 1 -c 0.9`) until
962 convergence, requiring 2, 3, and 3 rounds for the human-complete, other-complete, and human-partial
963 sets, respectively. This step yielded 0.154B, 7.4B, and 10.8B clusters, each of which with its representative
964 sequence, termed Logan90.

965 Next, the Logan90 databases were further clustered at 50% identity with the same coverage criterion
966 until convergence, requiring 1, 3, and 2 rounds for the human-complete, other-complete, and human-partial
967 sets, respectively. To ensure that dividing sequences into chunks during the final 50%-identity clustering
968 rounds (other-complete rounds 2–3; human-partial rounds 1–2) did not prevent the detection of sequences
969 belonging to the same cluster, we shuffled sequences between chunks. The final representative sets, termed
970 Logan50 databases, contained 0.07B, 3B, and 1.8B sequences, corresponding to overall reductions of 99.8%,
971 97.3%, and 99.4%. Altogether, we reduced 445.3B initial sequences to 4.87B representative sequences, a
972 total reduction of 98.9%.

973 3.13.1 Logan50 clustering enhances protein structure prediction by increasing multiple 974 sequence alignments (MSAs) diversity

975 To demonstrate the utility of Logan50 for downstream applications, we tested whether they can improve
976 protein structure predictions with AlphaFold2/ColabFold [33], focusing on viral proteins whose structures
977 are hard to predict with the default ColabFold DB. The Big Fantastic Virus Database [75] (BFVD)
978 comprises over 351,000 viral protein structures, generated by augmenting ColabFold default MSAs with
979 homologous proteins identified in Logan's contigs. In the BFVD paper, identification of homologs was
980 carried out over the entire 0.9 petabytes of Logan contigs. Here, we tested whether the reduced Logan50
981 dataset could achieve comparable performance while reducing the search space by approximately 3,500-
982 fold compared to using the entire Logan corpus.

983 We obtained 100 viral proteins from BFVD, previously characterized by poor-quality MSAs using the
984 default ColabFold DB. For each protein, we generated MSAs using MMseqs2 (`mmseqs search -a -s 8.5`
985 `--num-iterations 2 --max-seqs 1000`) against the other-complete Logan50 database, and compared
986 them to the default ColabFold MSAs in terms of MSA quality, measured by the number of effective
987 sequences (Neff) with `hhmake` (Fig. 5b, left) and structural prediction (Fig. 5b, right) quality, measured
988 by the pLDDT metric. Protein structure prediction using either default MSAs or Logan50 MSAs was
989 done by AlphaFold2/ColabFold (`colabfold_batch --num-models 1 --model-order 3`).

990 We observed substantial improvements in both metrics. Average Neff increased from 2.19 (baseline)
991 to 4.89 (Logan50 MSAs), while mean pLDDT scores improved from 46.7 ("very low") to 88.6 ("high").
992 Combining ColabFold and Logan50 MSAs further increased these values to a Neff of 5.18 and mean

993 pLDDT of 89.02, approaching the improvement achieved with the full Logan corpus plus ColabFold
994 sequences in the original BFVD paper (mean Neff = 4.39, mean pLDDT = 92.88). Thus, on par significant
995 gains in structural prediction can be achieved efficiently using the thousand times smaller Logan50, with
996 substantial reductions in computational requirements.

997 3.14 Expanding Obelisk Diversity with Logan

998 3.14.1 Reconstructing the Original Obelisk Database

999 A total of 1,744 Obelisk clusters (80% nucleotide clustering threshold), representing 7,202 circular genomes,
1000 were taken from the initial Obelisk study [9]. These 1,744 centroid were screened for false positives using
1001 several sequence alignment, structural homology and HMM verification steps (Data Availability, Obelisk
1002 Data Methods). The final clustering analysis resulted in 1,284 backward-compatible 60% Oblin-1 clusters
1003 and 1,965 90% clusters. Each of the 1,284 60% centroids also functioned as a 90% centroid, maintaining
1004 consistency in the database. This re-annotated database is designated as ‘Obelisk DB Legacy’.

1005 3.14.2 First Obelisk Expansion

1006 The Obelisk DB Legacy 1,744 Oblin-1 centroid sequences were used as query sequences for alignment to
1007 all Logan contigs using minimap2 (Section 3.2). The results of the Logan search were retrieved and split
1008 into circular and non-circular sequences by screening for 30-mer repeats at the ends of the Logan contigs.
1009 For circular sequences, the 30-mer repeat was trimmed, yielding 12,690 circular and 312,728 non-circular
1010 contigs.

1011 Concurrently, the 7,195 Oblin-1 sequences from ‘Obelisk DB Legacy’ were clustered at 95% identity
1012 using USEARCH cluster_fast, resulting in 2,170 clusters. Using Clustal Omega [76] and HMMER [63],
1013 a Hidden Markov Model was constructed from this clustered dataset. The circular sequences from the
1014 Logan output underwent ORF calling using EMBOSS getorf with circular genome parameters (getorf
1015 -circular Yes), and the resulting ORFs were incorporated into the ‘Obelisk DB Legacy’ HMM model
1016 through an iterative search and alignment protocol (Data Availability, Obelisk Data Methods).

1017 This procedure resulted in a comprehensive Obelisk database containing 117,838 Oblin-1 proteins,
1018 along with a high-quality MSA comprised of 45,170 Oblin-1 sequence alignments. All 117,838 Oblin-1
1019 proteins were verified to contain the Domain A motif, and were aligned against the BLAST nr database
1020 (May 2024) using DIAMOND BLASTp [11] with parameters –masking 0 –unal 1 –sensitive -c1 -k1 -b 5
1021 –threads 16. This analysis confirmed that no Obelisk sequences produced significant hits in the BLAST
1022 nr database, indicating that all sequences represent completely novel genetic elements.

1023 3.14.3 Construction of Obelisk DB v1

1024 The original Obelisk sequences from ‘Obelisk DB Legacy’ were excluded from the 117,838 Obelisk se-
1025 quences, yielding 110,643 novel sequences. These novel sequences were aligned to the 1,284 60% centroid
1026 sequences from the Legacy Database using USEARCH usearch_global with 60% identity threshold (use-
1027 arch -usearch_global -id 0.6). Sequences that aligned within known 60% clusters were subsequently aligned
1028 to the corresponding 90% centroid sequences within their respective 60% clusters using USEARCH use-
1029 arch_global with 90% identity threshold (usearch -usearch_global -id 0.9).

1030 This two-step alignment procedure classified all sequences into three distinct categories: (1) members
1031 of known 60% and 90% clusters, (2) members of known 60% clusters that represent novel 90% clusters,
1032 or (3) members of novel 60% clusters. Using the alignment identities and closest matches identified by
1033 USEARCH, sequences were assigned to their appropriate clusters.

1034 For novel 60% or 90% clusters, centroid selection prioritized circular sequences, followed by sequence
1035 length as the secondary criterion. In clusters lacking circular sequences, the longest sequence was des-
1036 ignated as the representative centroid. Metadata, including amino acid sequence, nucleotide sequence,

1037 circularity status, SRA accession numbers, BioProject information, and additional relevant data, was
1038 appended to each database entry. The resulting database was designated ‘Obelisk DB Logan v1’.

1039 3.14.4 Database Refinement with Logan v1.1 Contigs

1040 As Logan contigs were updated from v1 to v1.1, we also updated the Obelisk database. Centroids from
1041 the 90% Obelisk clusters were extracted and aligned to Logan v1.1 contigs using minimap2 (Section 3.2).
1042 The results were processed using the same methodology described in the previous section. Circular and
1043 non-circular sequences were separated, and k -mer repeats were removed following the established protocol.

1044 The v1.1 Obelisk database was constructed as follows. Beginning with the original 7,195 sequences and
1045 the previously established high-quality 45,000-sequence MSA, the same iterative alignment methodology
1046 with tapered e-value cutoffs was executed. Circular sequences were processed first, followed by non-circular
1047 sequences. In total 67,454 Oblin-1 sequences were captured.

1048 BLAST alignment analysis revealed no significant hits for all proteins except one sequence, which
1049 showed a very distant, low e-value alignment. The constructed HMM model was used to query all 67,454
1050 sequences, and only 550 Oblin-1 sequences exhibited alignment e-values greater than e-5 to the model;
1051 these sequences were retained to preserve diversity. This database was designated ‘Obelisk DB Logan
1052 v1.1’.

1053 Comparative analysis using DIAMOND BLASTp between versions 1 and 1.1 confirmed that all se-
1054 quences present in version 1 were accounted for in version 1.1, despite the total number of Obelisk
1055 sequences being approximately halved. This reduction can be attributed to the differences in average se-
1056 quence lengths between the two databases. Version 1 had an average amino acid length of 136.7 residues,
1057 while version 1.1 had an average length of 186 residues. This increase in sequence length explains the
1058 decrease in individual sequence counts, while maintaining comprehensive coverage.

1059 The same methodology described in ‘Construction of Obelisk DB v1’ was applied and the complete
1060 ‘Obelisk DB v1.1’ was assembled with all associated metadata incorporated.

1061 3.15 Genomic Expansion of P4 Phage Satellites

1062 We first gathered all RefSeq protein sequences for each of the seven core components in a P4 phage satel-
1063 lite [34] (Counts: alpha: 1911, ash: 1995, alpA: 2008, Sid: 2041, Delta: 2089, Psu: 1891, Integrase: 2097).
1064 Sequences were clustered for each core component individually via UCLUST v11.0.667 `cluster_fast` [60]
1065 with identity threshold set to 0.9, then cluster representatives from each of the core components were
1066 aligned against Logan v1.1 contigs via DIAMOND BLASTX, as per Section 3.2, with sensitive mode
1067 and e-value minimum of 1e-8, for a total of 107.6 million putative P4 satellite contigs. Contigs were
1068 next refined with HMMER v3.4 [63], by 6-frame translating each of the putative P4 contigs with seqkit
1069 v0.10.0 `translate` [49], then running HMMER `hmmscan` on each of the translated sequence against
1070 each of the core component protein databases. A total of 210,895 contigs contained a hmmscan hit
1071 for each core component with a minimum e-value of 1e-5 and a coverage of at least 40% of the database
1072 sequence. Finally, these contigs were clustered with MMseq2 v15 [77] with minimum identity threshold
1073 of 0.999 and coverage of 0.8, curating 16,215 cluster representatives. Proteins were detected within each
1074 representative with Prodigal [72] (version 2.6.3) and P4 phage satellites were detected and characterized
1075 with SatelliteFinder v1 [34] with default parameters.

1076 For the pangenome curve, proteins from all RefSeq P4 phage satellites of types A, B, or C were clustered
1077 with MMseqs2. Clustering was then repeated with each cluster representative and proteins from Type A
1078 and B Logan contig containing less than 30 genes (15,653 contigs). Any cluster containing only Logan
1079 contigs was considered a new gene family. For the whole-genome reciprocal relatedness (wGRR) analysis,
1080 full proteomes from the 16,653 Logan P4 contigs and the 3,437 RefSeq P4 genomes were extracted by
1081 detecting the first and last protein in the P4 satellite region. A wGRR matrix was formed by calculating
1082 the fraction of bi-directional best hits weighted by the sequence identity for each genome pair in an all-
1083 vs-all fashion. Hierarchical clustering was performed on the matrix and the corresponding heatmap was

1084 rendered with with seaborn's `clustermap` function, using the Ward clustering algorithm. All scripts can
1085 be accessed at <https://github.com/kdcurry/P4-logan>.

1086 3.16 Plasmid identification and clustering

1087 We identified 3,885,511 circular contigs ($\geq 1\%$ kb) from version 1.0 assembly graphs of 252,507 samples,
1088 including bacterial, archaeal, and metagenomic datasets (code available at https://gitlab.pasteur.fr/rchikhi_pasteur/logan-circles). Contigs were processed with tr-trimmer (version 0.1.0, parameters:
1089 `-c -x -l 31`) to discard sequences with low-complexity repeats spanning $> 50\%$ of terminal 31-bp re-
1090 peats and to trim these repeats from the 3' ends. To mitigate gene truncation, sequence breakpoints were
1091 shifted to intergenic regions (code available at <https://github.com/apcamargo/reorient-circular-seq>). Plasmids were then identified using geNomad [78] (version 1.8.1, database version 1.7, `end-to-end`
1092 command, parameters: `--enable-score-calibration --max-fdr 0.01`). To minimize false positives,
1093 we only kept plasmids that met two criteria: (1) encode at least one protein matching a HMM from a
1094 curated set of 193 plasmid hallmark protein profiles (see Data Availability); (2) encode no more than
1095 two proteins matching HMMs of near-universal single-copy orthologs from the BUSCO v5 [79] odb10 [80]
1096 datasets of *Bacteria* and *Archaea*. Gene prediction was carried out using pyrodigal-gv [81, 78] (ver-
1097 sion 0.3.2), and protein sequence matching to HMMs was performed using PyHMMER's [82] (version
1098 0.10.15) `hmmssearch` function, applying gathering cutoffs to the HMMs of plasmid hallmarks (param-
1099 eter: `bit_cutoffs="gathering"`) and BUSCO bitscore cutoffs to the HMMs of near-universal single-copy
1100 orthologs. Identified plasmids were assigned to replicon families using MOB-typer [83] (version 3.1.9).
1101 AMR genes were annotated using AMRFinderPlus [84] (version 4.0.3), and AMPs were identified with
1102 Macrel [85] (version 1.5.0).

1103 To cluster plasmids, we first computed pairwise sequence similarities using BLAST [57] (version 2.16.0,
1104 parameters: `-task megablast -evalue 1e-5`) and the `anicalc.py` script from CheckV [86] to compute
1105 pairwise similarity metrics. We then constructed a similarity graph connecting pairs of plasmids exhibiting
1106 sequence identity $\geq 90\%$ and bidirectional alignment coverage $\geq 90\%$, and clustered the plasmids using
1107 the pyLeiden [87] tool.

1108 We evaluated Faith's phylogenetic diversity of selected replicases and relaxases (Data Availability,
1109 Plasmid PD Table) from newly identified plasmids and complete plasmids retrieved from PLSDB (release
1110 2024.05.31_v2) and IMG/PR. Genes encoding these proteins were identified using `hmmssearch`, and multi-
1111 ple sequence alignments were generated with PyHMMER's `hmmlalign` function (parameters: `trim=True`,
1112 `all_consensus_cols=False`). These alignments were then trimmed with the `gappyout` algorithm from
1113 PytrimAl [88] (version 0.8.0) and phylogenetic trees were inferred with FastTree [89] (version 2.1.11).

1114 We surveyed plasmid presence across all samples by mapping contigs from version 1.1 assemblies to
1115 plasmid cluster representatives using minimap2 (version 2.28, parameters: `-x sr --sam-hit-only -a`).
1116 A plasmid cluster was considered present in a sample if $\geq 75\%$ of its length was covered by alignments
1117 from that sample.

1118 3.17 Antimicrobial resistance genes discovery and analyses

1119 We analysed the presence of antimicrobial resistance (AMR) genes across 26.7 million SRA accessions
1120 via the Logan v1.1 contigs. AMR gene hits were identified by aligning the CARD nucleotide database
1121 (version 3.3.0) to Logan contigs using minimap2, as per Section 3.2 (see Data Availability). Alignments
1122 were filtered to contain only those sequences with >100 bp length and $>80\%$ identity. SRA meta-
1123 data and extended geolocation data (see Data Availability) were used to classify information on CARD
1124 alignment SRA accession hits. For analyses based on organism classification, datasets were classified as
1125 organism type Metagenome if SRA metadata field `organism` contains the string "metagenome", or if the
1126 field `librarysource` contains "METAGENOMIC" or "METATRANSCRIPTOMIC"; otherwise, acces-
1127 sions were classified as Isolate. Metagenome categories were classified according to the `organism` and
1128 `librarysource` fields, dividing it into the 6 top categories: human, livestock, soil, marine, freshwater,

wastewater. Plasmid metadata was extended analogously to SRA samples. For Extended Data Fig. 7d-g), the SRA-CARD alignment dataset was filtered to include only samples with known collection dates, geolocation, metagenomic origin, and those more likely to contain whole genome/transcriptome assay types (WGS, WGA, RNA-Seq, etc.). All code and datasets can be accessed in: https://github.com/mmontonerin/logan_AMR

3.18 Comparison of Logan metagenomes with GenBank WGS via FracMinHash sketches

FracMinHash sketches were created with `sourmash` [90] for each Logan unitigs accession whose SRA metadata information indicated it was a metagenome (total of 4,792,069 accessions). We used a k -mer size of 31 and a scale factor of 1,000. The resulting sketches, comprised of 449,087,511,713 hashes with duplicates, were then placed in a duckDB database [91]. Code to process these data are available at: https://github.com/KoslickiLab/ingest_logan_yacht_data.

We then crawled the GenBank Whole Genome Shotgun (WGS) FTP server, downloaded and sketched each `*.fsa_nt.gz` assembly with `sourmash` using the same k -mer size of 31 and scale factor of 1,000. Of the 2,055,047 `*.fsa_nt.gz` files discovered on this FTP server, 5 led to HTTP error code 404 when attempting to download them. Signatures were stored as compressed `*.sig.zip` archives. We extracted all 64-bit FracMinHash hash values, a total of 32,701,966,322 hashes with duplicates, and partitioned them by a low-bits bucket function to enable external-memory de-duplication. We then reduced each bucket to its exact set of 8,679,649,739 unique hashes and wrote a partitioned Parquet dataset which we ingest into DuckDB. To compare against the Logan metagenome sketches, we attached the Logan metagenome FracMinHash sketches, and again bucketed into Parquet files and computed set differences with bucket-wise anti-joins in DuckDB to report the number of 31-mers (appearing twice or more) unique to each dataset: 32,287,730,882 in the Logan metagenomes not in GenBank WGS, and 7,020,467,844 in GenBank WGS not in the Logan metagenome sketches. Since a scale factor of 1,000 was used to form the sketches, multiplying these number of hashes by 1,000 results in the estimated total number of 31-mers, each appearing twice or more in each dataset. All code can be accessed at: https://github.com/KoslickiLab/GenBank_WGS_analysis.

3.19 Geographic Metadata Extraction, Inference and Enrichment

99.7% of SRA submissions are associated with a BioSample <https://www.ncbi.nlm.nih.gov/biosample> record. BioSamples contain zero or many attribute [name, value] tuples with submitter-supplied metadata that describes the biological sample. A full XML dump of the BioSample database was retrieved on June 28, 2024, comprising 39,448,576 records with 568,885,433 attribute entries. A subset of attribute names likely to contain geographical information was extracted using a large language model (gpt-3.5-turbo-0125) combined with manual curation of the most frequently occurring ones. The distribution of attribute names is heavily skewed and long-tailed: the eight most common names account for 92.84% of all entries identified as containing geographical information.

Attribute values corresponding to this subset of attribute names were then used to infer geographic coordinates. A deep learning classifier partitioned the values into three distinct categories: values likely to contain numerical latitude and longitude pairs (coordinates); values likely to reference a location by its common name (place names); and values which were explicitly annotated but whose meaning is uninformative, e.g. “N/A”, “null”, “undefined”, etc. Coordinates were resolved directly to points on Earth under the WGS 84 coordinate system. Place names were converted to coordinates using three different geolocation services (Azure Maps <https://azure.microsoft.com/en-us/products/azure-maps>; and Esri, HERE through AWS Location <https://aws.amazon.com/location/>). To quantify confidence in these derived coordinates, a score of 0 to 6 was assigned depending on whether the locations returned by the three services were within 8 km of each other (up to 3 points), and whether they were found within the political boundaries of the same country (up to 3 points). Values classified as uninformative were discarded. In total, geographic coordinates were obtained for 26,962,465 (68.34%) of BioSample entries.

1178 The geographic dataset was further enriched by cross-referencing it with the following publicly available
1179 resources: *ASTER Global Digital Elevation Model* <https://cmr.earthdata.nasa.gov/search/concepts/C1711961296-LPCLOUD.html> to extract elevation in meters above mean sea level; *World Administrative Boundaries* <https://public.opendatasoft.com/explore/dataset/world-administrative-boundaries/export/> to assign political boundaries such as countries and regions; and *WWF Terrestrial Ecoregions of the World* <https://www.worldwildlife.org/publications/terrestrial-ecoregions-of-the-world> to classify BioSamples into 14 distinct biomes characterized by their unique biodiversity and environmental conditions.

1186 4 Data Availability

1187 Logan is publicly available [37] on AWS Registry of Open Data at <https://registry.opendata.aws/pasteur-logan/>. There are no egress charges and anonymous access is permitted. A data access tutorial
1188 is provided at <https://github.com/IndexThePlanet/Logan>. PETadex data can be found at <https://github.com/ababaian/petadex>. AMR datasets can be found at <s3://logan-pub/paper/AMR>. Obelisk
1189 data and code can be found at <s3://logan-pub/paper/Obelisk>, and the Obelisk Data Methods is
1190 methods.pdf within this folder. Plasmid sequence and metadata can be found at <s3://logan-pub/paper/plasmids>, and the Plasmid PD Table is plasmid_phylogenetic_diversity.xlsx within this
1191 folder.

1195 5 Acknowledgments

1196 We are grateful to the entire team managing the NCBI SRA / EBI ENA and the biology community for
1197 data sharing. We thank Dorian Schaal, Adrien Lainé, Coral Kennett, Candi Jeronimo, Dave Maurer, and
1198 Morgan Lim from Amazon Web Services (AWS) for support. Thomas Menard, Stéphane Fournier and
1199 the HPC Core Facility from Institut Pasteur for IT support. Peter Schmiedeskamp, Chris Stoner, Erin
1200 Chu and Beryl Rabindran from AWS Registry of Open Data for data hosting. Ryan Connor and Yuryi
1201 Skripchenko from NCBI for assistance with SRA operations. Karen McGregor, Jerry Morey and Venkat
1202 Malladi from Microsoft for Azure support. Matthieu Falce for custom AWS tooling. Stephen Nayfach
1203 for feedback on the analysis of plasmid genomes. Karin Steffen and Zamin Iqbal for feedback on the
1204 manuscript. Administrative support was provided by Mélanie Ridel, Loïc Orellou, and Florence Percie
1205 du Sert.

1206 6 Funding

1207 Computing resources were provided by the University of Toronto Cloud Research Lab at The Donnelly,
1208 powered by AWS. R.C. was supported by ANR grants ANR-22-CE45-0007, ANR-19-CE45-0008,
1209 PIA/ANR16-CONV-0005, ANR-19-P3IA-0001, ANR-21-CE46-0012-03, and Horizon Europe grants No.
1210 872539, 956229, 101047160 and 101088572 (ERC IndexThePlanet, supporting also K.D.C. and R.F.).
1211 A.B was supported by Canadian Institutes for Health Research (CIHR) project grant PTJ-496709, and
1212 as a Canadian Institute for Advanced Research (CIFAR) Global Azrieli Scholar of the CIFAR Fungal
1213 Kingdom: Threats & Opportunities program. C.A.L. is supported by National Institutes of Health
1214 grants R00HG012579 and P30CA008748 as well as a Michelson Medical Research Foundation Award.
1215 A.P.C., C.J.M. and M.B.F. were supported by the US Department of Energy Joint Genome Institute
1216 (<https://ror.org/04xm1d337>), Office of Science user facilities, operated under contract no. DE-AC02-
1217 05CH11231, and Office of Biological and Environmental Research (BER) as part of BER's Genomic
1218 Sciences Program (GSP) under FWP 70880. P.P. is supported by Inria challenge program OmicFinder
1219 and ANR-19-CE45-0008. G.W.B. was supported by the Natural Sciences and Engineering Research Council
1220 of Canada (RGPIN-2017-06855) and a Tier 1 Canada Research Chair (CRC). T.L. was supported by

1221 ANR-19-CE45-0008. M.M-N. was supported by Leverhulme Centre for the Holobiont. J.S. was supported
1222 by the Natural Sciences and Engineering Research Council of Canada (Canada Graduate Scholarship -
1223 Master's). D.K. was supported by NIH NIGMS grant R01GM146462. P.M. was supported by NSF
1224 grants DBI2138585 and OAC1931531, and NIGMS/NIH grant R01GM146462. J.A.M.S was supported
1225 by ANR-23-CE20-0046 01 TRIADE. D.P.A. was supported by NIH grant U19AI144297. K.S. is funded by
1226 the Canadian Foundation for Innovation (CFI) and the Canadian Institutes of Health Research (CIHR,
1227 grant number 186156). P.J.R. is funded via CIHR (grant numbers 186156 and 197950) and is a CRC in
1228 Chemical Genetics.

1229 **7 Contributions**

1230 R.C. and A.B. conceived and led the study. R.C., G.A., M.H., A.K. and B.R. designed and implemented
1231 the Logan assemblage cloud infrastructure. A.K. developed the f2sz software. R.C. designed and im-
1232 plemented the Logan contigs mining analyses. T.L. and P.P developed Logan-Search. A.B., R.L-K.,
1233 R.C.E, J.S., R.C., K.S., P.J.R. and G.W.B analyzed the plastic-active enzymes. C.A.L. analyzed the
1234 viral reactivation. M.S. and J.L. clustered the Logan proteins, and M.S. analyzed the viral proteins case
1235 study. R.F. created protein embeddings. P.G. analyzed the Obelisks. K.D.C., J.A.M.S and E.P.C.R.
1236 analyzed the P4 satellites. A.P.C., C.J.M. and M.B.F. analyzed the plasmids. M.M-N., D.P.A. and S.M.
1237 designed the AMR study, M.M-N. analyzed the AMR in Logan contigs, M.M-N. and A.P.C. analyzed the
1238 AMR in plasmids. D.K. analyzed Logan metagenome contigs and WGS sketches. A.M-T. created the
1239 SRA geographic metadata dataset. R.C., A.B., R.C.E., C.A.L, P.P., M.S., R.L-K., A.P.C., M.M-N., P.G.,
1240 K.D.C., E.P.C, D.K., P.M., and A.M-T. wrote the manuscript.

1241 All authors contributed to, and approved the manuscript.

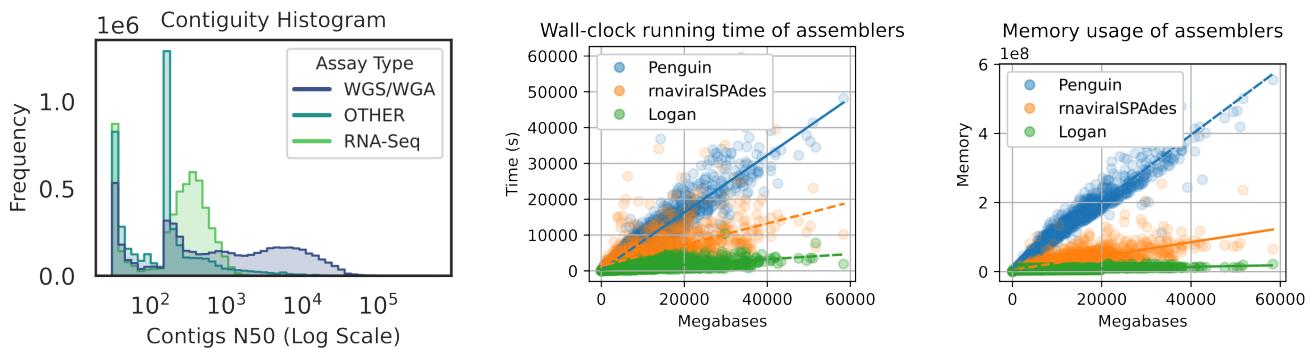
1242 **8 Competing interests**

1243 The authors declare no competing interests.

1244 **9 Materials & Correspondence**

1245 Correspondence and requests for materials should be addressed to Rayan Chikhi or Artem Babaian.

1246 Extended Data



Infrastructure statistics for computation over the entire SRA

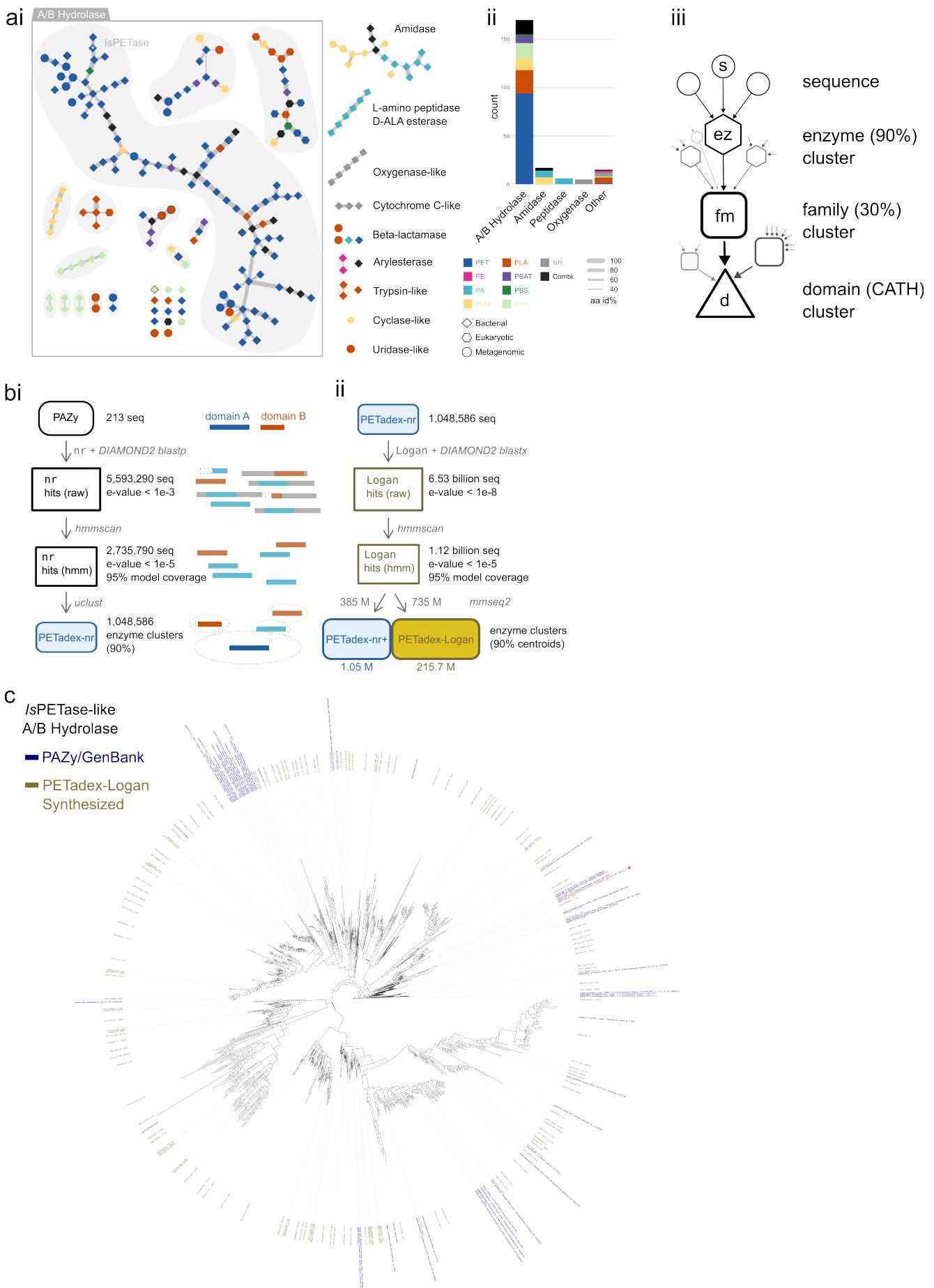
Global statistics		Run 6 statistics	
Input SRA Accessions	27 million	Input data	19.6 petabases
Input SRA size	50 petabases	Runtime*	7 hours
Total CPU Hours	~30 million	Peak Number of Instances	73,100
Number of Runs	6	Peak Number of vCPUs	2.18 million
Total Runtime	30 hours	Peak Total EBS storage	52 petabytes

Extended Data Fig. 1: **Logan assembly performance and computational statistics for processing the entire SRA.** This figure details the performance benchmarks of the Logan pipeline and quantifies the cloud computing resources used to assemble 27 million SRA datasets.

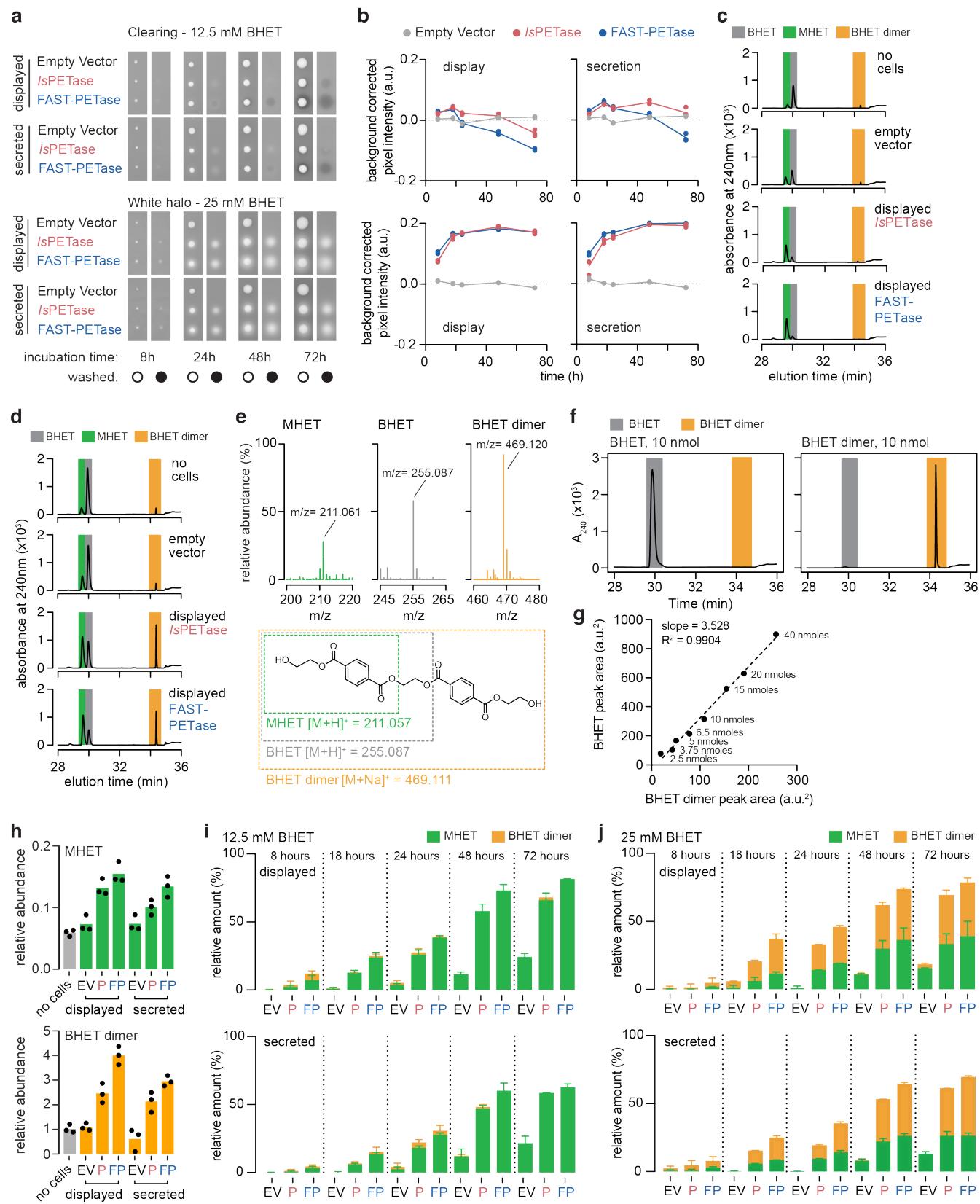
(Top Left) A histogram showing the distribution of assembly contiguity, measured by contig N50, across all Logan assemblies. Assemblies are categorized by input SRA assay type, showing that Whole Genome Shotgun (WGS/WGA) samples generally produce more contiguous assemblies than RNA-Seq or other samples, as expected.

(Top Middle and Right) Performance benchmarks comparing the Logan assembly pipeline to other state-of-the-art short read metagenome assembly tools (Penguin, maviralSPAdes). Logan pipeline demonstrates significantly lower wall-clock running time (middle) and memory usage (right) across a range of input data sizes, highlighting its efficiency.

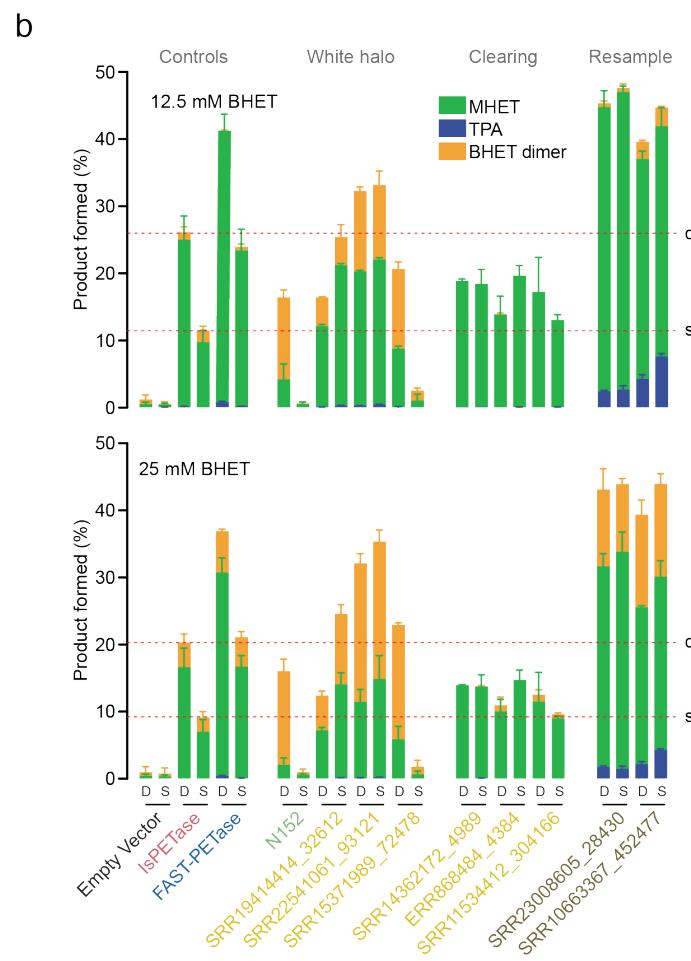
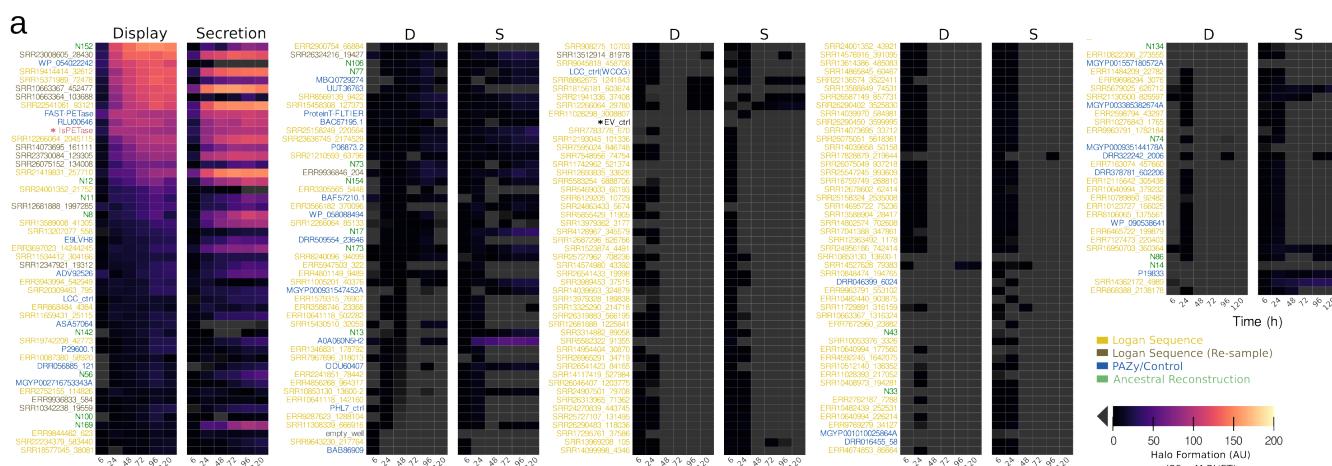
(Bottom) Statistics from the full-scale production run. Global statistics summarize the total compute effort, including processing 50 petabases of input data over 30 million CPU hours. The vCPU Usage Over Time plot for the main production run illustrates the dynamic allocation of cloud processors, peaking at over 2.18 million vCPUs. Run 6 statistics detail the single largest run, where 19.6 petabases of data were assembled in just 7 hours.



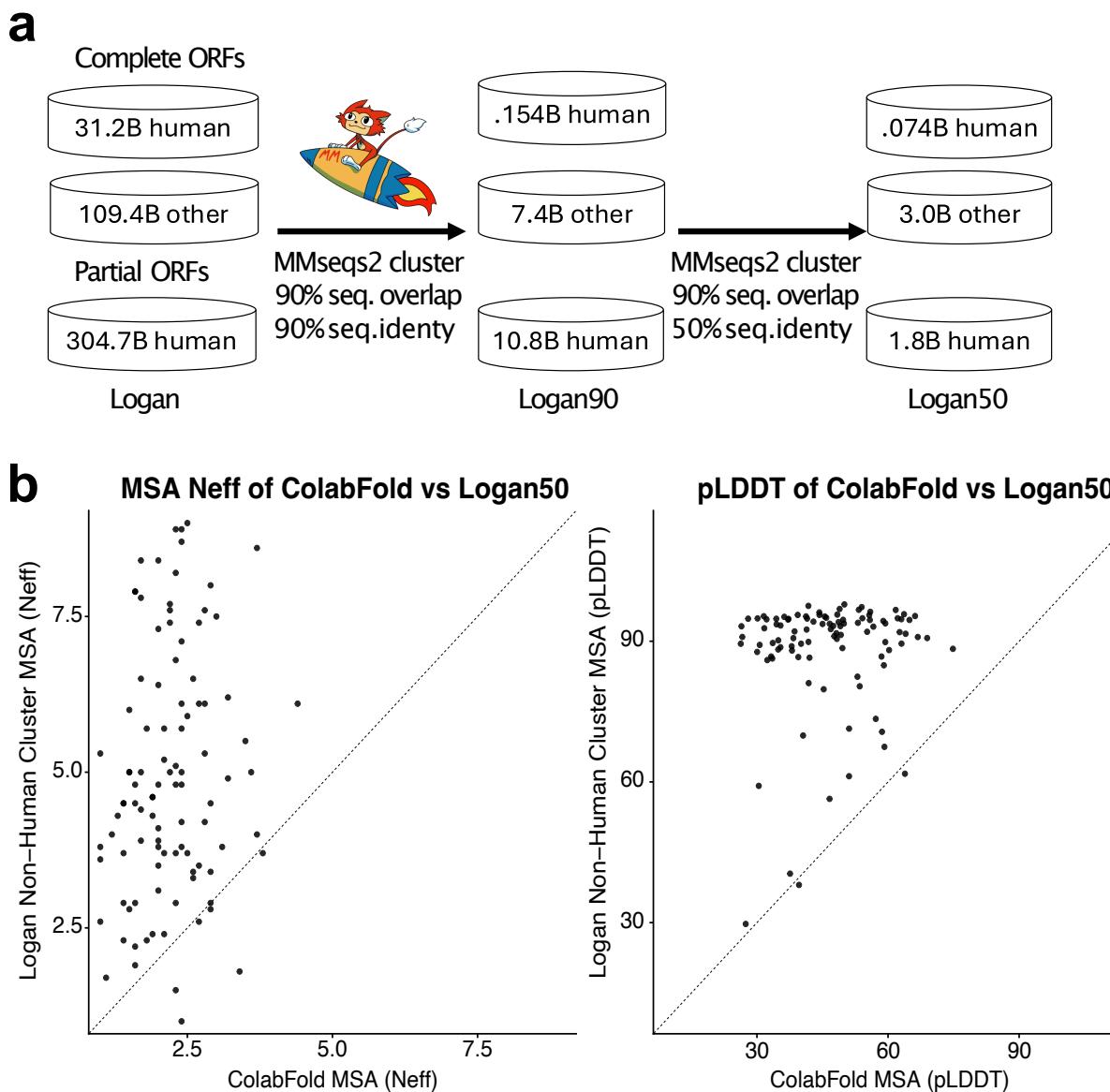
Extended Data Fig. 2: **PETadex-Logan Workflow.** **(a)** Characterization of the initial 213 plastic-active enzymes from the PAZy database. **(i)** A network graph showing sequence similarity between the known enzymes, colored by protein family. **(ii)** Bar chart showing the distribution of these enzymes across protein families and the types of plastics they degrade. **(iii)** Schematic of the hierarchical clustering strategy used to group sequences at the enzyme (90% identity), family (30% identity), and domain (CATH) levels. **(b)** The two-stage deep homology search pipeline. **(i)** The first stage queried PAZy sequences against the NCBI nr database. After filtering for domain integrity and clustering, this step yielded 1.05 million enzyme clusters, creating the PETadex-nr dataset. **(ii)** In the second stage, PETadex-nr was queried against the entire Logan assembled contigs, identifying 735 million novel sequences and massively expanding the diversity into the final PETadex-Logan dataset. **(c)** A phylogenetic tree of the *IsPETase*-like A/B Hydrolase clade. The tree visually demonstrates the expansion of sequence diversity uncovered by the Logan search compared to the previously known diversity from public databases like PAZy and GenBank (blue labels). Sequences selected for experimental evaluation are labeled.



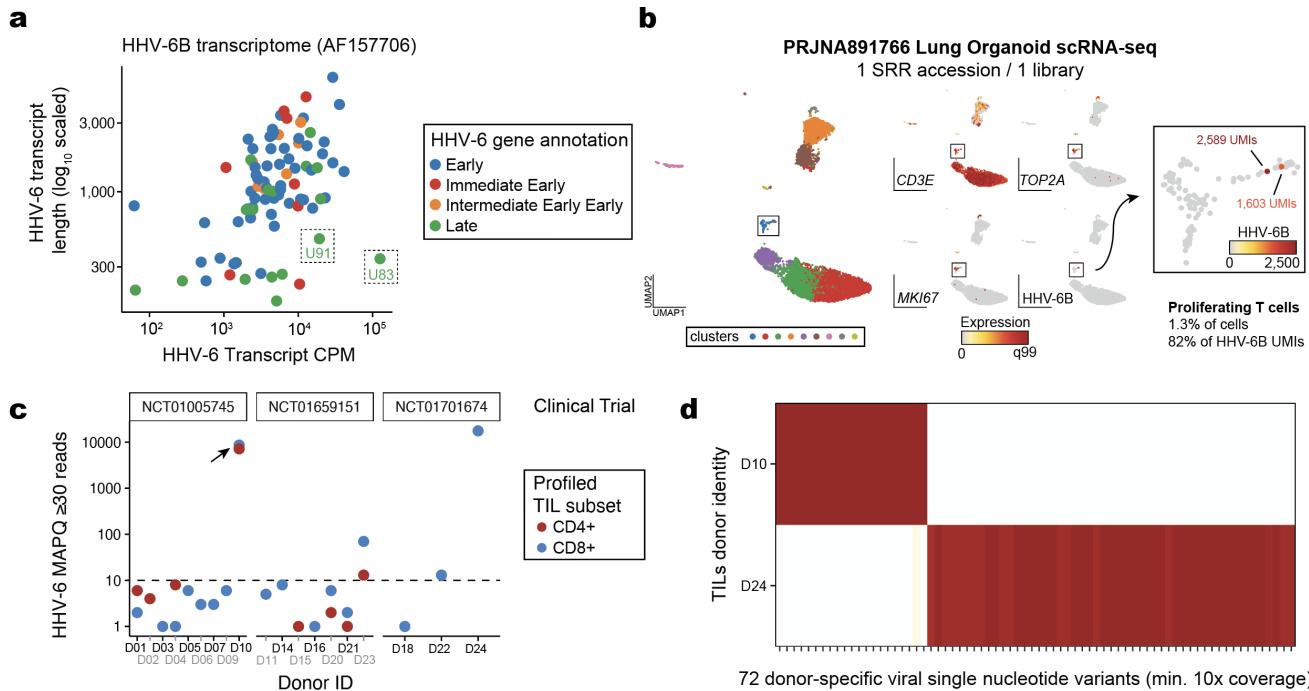
Extended Data Fig. 3: **PETase Halo Assay.** **(a)** Clearing and white halo detection with yeast colonies expressing surface displayed or secreted PETase enzymes. Yeast strains were robotically pinned onto YPD medium containing 12.5 mM or 25 mM BHET and incubated for 8 to 72 hours at 30 degrees Celsius. Plates were imaged before (open circle) and after (closed circle) washing colonies off the plate. **(b)** Clearing and white halo quantification from washed YPD plates containing 12.5 mM or 25 mM BHET. Pixel intensity of the colony area was measured using a custom R pipeline (see Methods) on images from (a). Data is depicted as the background normalized median pixel intensity under each colony over time for the indicated PETases; n=3. **(c)** High-Performance Liquid Chromatography (HPLC) analysis of the clearing zone identifies the MHET reaction product. Agar plugs were excised from the plates in (a) for the displayed PETases , empty vector control, and regions with no yeast cells, after 72 hours of yeast growth on 12.5 mM BHET, and dissolved in DMSO prior to HPLC analysis. Representative chromatograms are shown; $n \geq 2$. Coloured shading indicates the identity of each peak. **(d)** HPLC analysis of the white halo identifies MHET and BHET dimer reaction products. Agar plugs were excised from the plates in (a) for the displayed PETases , empty vector control, and regions with no yeast cells, after 72 hours of yeast growth on 25 mM BHET, and dissolved in DMSO prior to HPLC analysis. Representative chromatograms are shown; $n \geq 2$. Coloured shading indicates the identity of each peak. **(e)** Mass spectrometric (MS) analysis of MHET, BHET, and BHET dimer purified from a white halo extracted under a yeast colony expressing surface-displayed *IsPETase* after 24 hours on YPD plus 25 mM BHET. Representative spectra are shown with the mass to charge ratio of the most abundant component indicated; n = 3. The chemical structures of BHET, MHET and putative BHET dimer are shown along with their predicted ionized mass. **(f)** HPLC analysis of 10 nmol of HPLC-purified BHET, and 10 nmol of HPLC-purified BHET dimer. Coloured shading indicates the identity of each peak. **(g)** Correlation plot of absorbance peak areas from HPLC analysis of the indicated amounts of purified BHET and BHET dimer. The linear regression line is plotted. **(h)** MS quantification of MHET and BHET dimer from agar plug extraction. Agar plugs were obtained from an area of YPD plus BHET 25 mM with no yeast colony (no cells) or under the yeast colonies (after wash) containing the indicated constructs after 24 hours of incubation and processed as in (d). Relative abundance is plotted, expressed as a ratio between spectral counts for MHET (top) or BHET dimer (bottom) relative to the spectral counts obtained for BHET. EV: empty vector; P: *IsPETase*; FP: FAST-PETase. **(i,j)** Quantification of BHET conversion in clearing zones and white halos over time. Halos from the indicated strains, timepoints and BHET concentrations were processed as described in (c,d) and analyzed by HPLC. Peak area for each analyte (BHET, MHET, BHET dimer) was measured and expressed as a percentage relative to the sum of the peak areas for BHET+MHET+BHET dimer. EV: empty vector; P: *IsPETase*; FP: FAST-PETase.



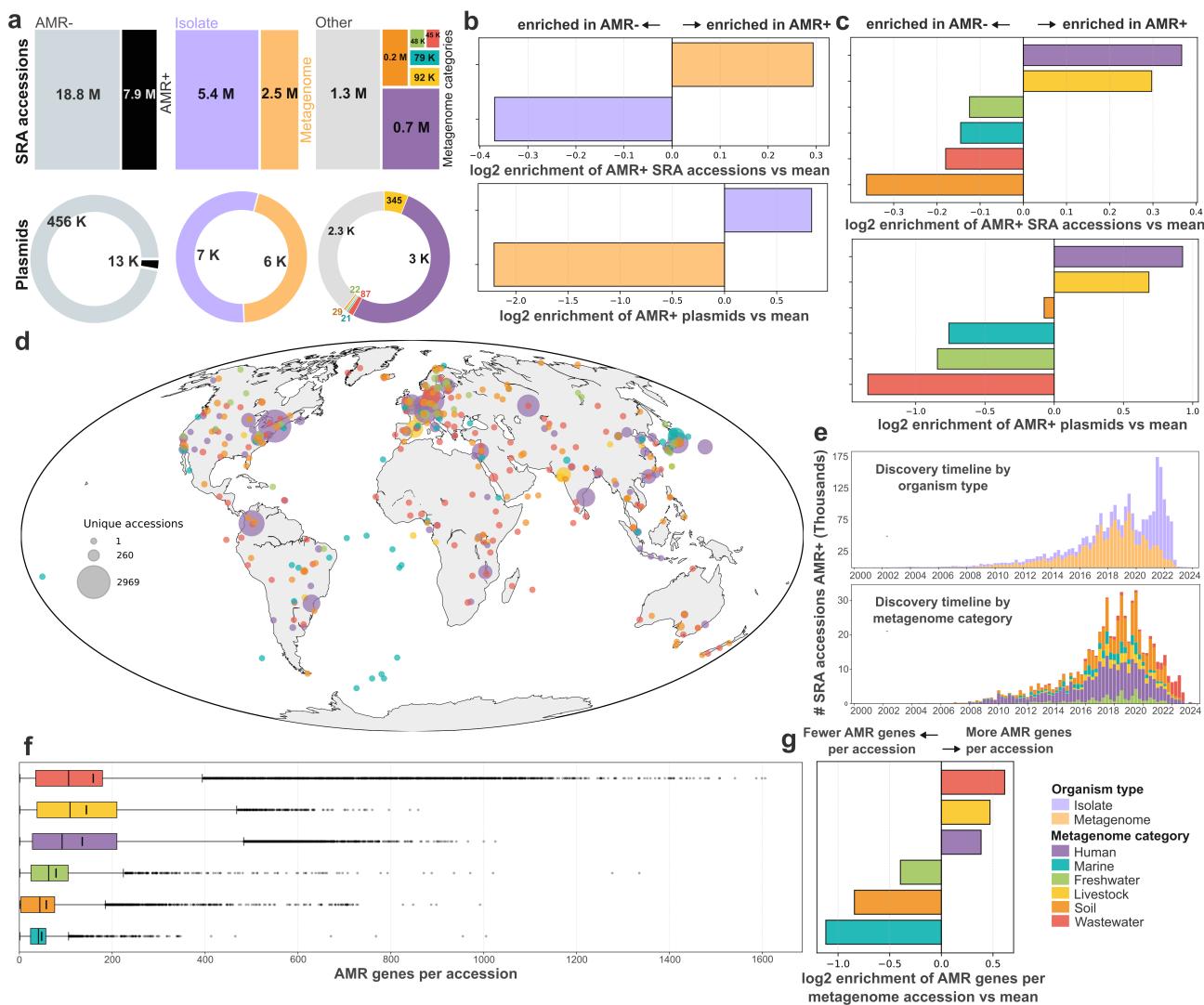
Extended Data Fig. 4: **High-throughput screening and HPLC validation of PETadex-Logan enzymes.** (a) Heatmap of the high-throughput activity screening results for Logan PETases and controls. Enzyme activity was measured as the background normalized median pixel intensity under each colony on YPD plates with 25mM BHET, at the indicated times, and in either surface displayed (D) or secreted (S) constructs. The heatmap shows the average of quadruplicate pixel intensity measurements (in arbitrary units, AU) after subtracting Empty Vector background values and scaling to approximately 100 units for *IsPETase* at 48 hours. This screen was used to identify the active candidates for quantitative analysis. (b) Quantification of BHET conversion in yeast strains expressing the top candidate PETase enzymes. Strains were grown to saturation in YPD medium prior to adding BHET at the indicated concentrations. BHET conversion reactions were allowed to proceed for 17 hours at 30°C, and culture supernatants were analyzed by HPLC. The peak area for each analyte (BHET, MHET, BHET dimer) was measured and expressed as a percentage of the sum of all peak areas normalized to 10^8 cells/ml, based on the cell concentration at the time of BHET addition. D: surface displayed enzyme; S: secreted enzyme.



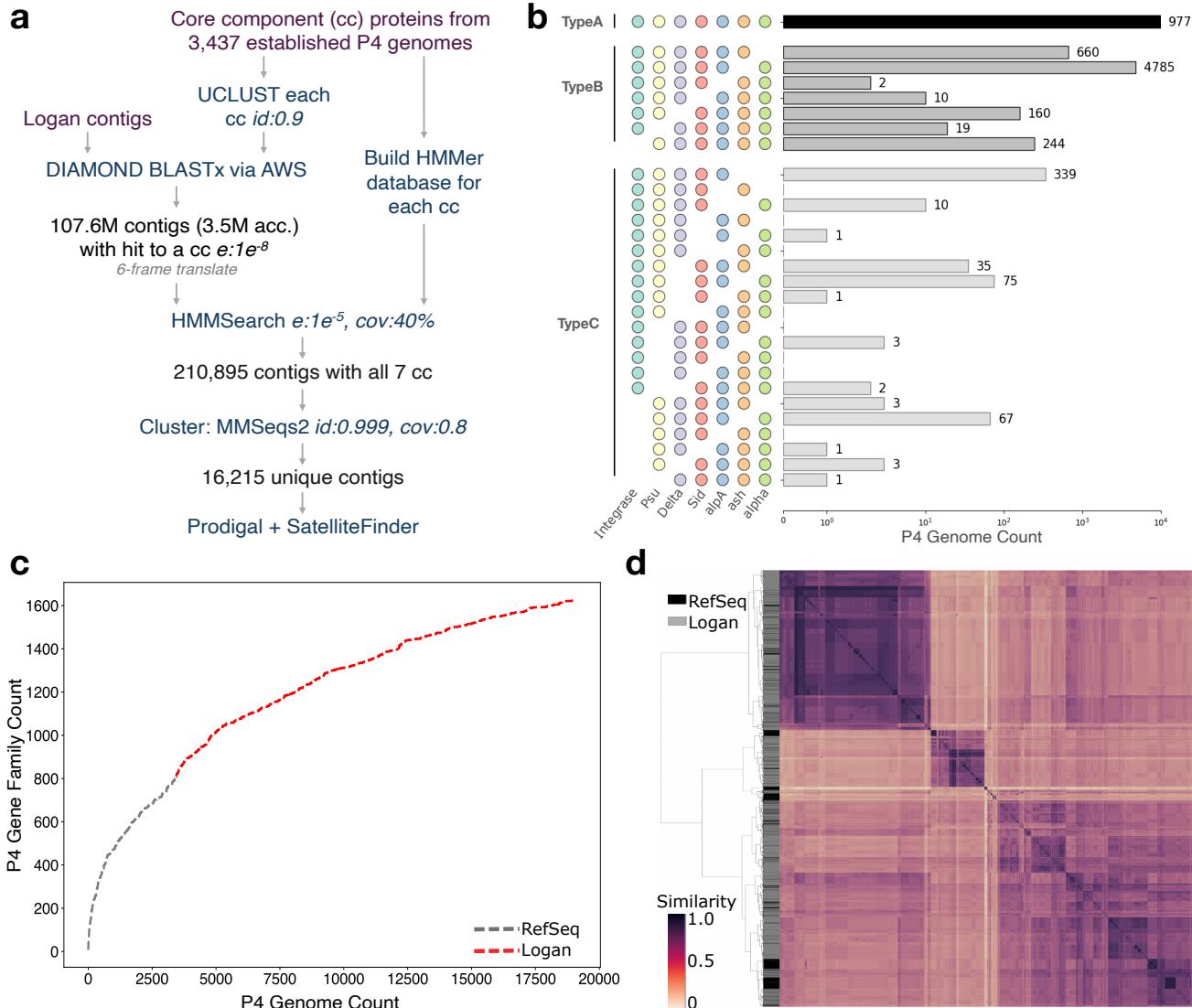
Extended Data Fig. 5: **Protein clustering workflow and its application in improving multiple sequence alignment (MSA) diversity.** (a) The workflow for creating the Logan90 and Logan50 clustered protein databases. Prodigal-predicted protein coding regions from all Logan contigs were first separated into 'human' and 'other' categories, based on SRA metadata associated with their contig and into 'complete' and 'partial', based on Prodigal's output. The proteins were then clustered using Linclust at 90% and subsequently 50% sequence identity to create representative protein sets for sensitive homology searches. The numbers indicate billions (B) of proteins at each stage of the workflow. (b) A case study demonstrating the value of Logan50 for enhancing MSA diversity (Neff, left panel) and improving structure prediction quality (pLDDT, right panel) of 100 viral proteins with low-quality MSAs from the default ColabFold database. We performed sensitive, iterative profile searches against the other-complete Logan50 database (y-axis) using MMseqs2 and compared the results to those from the default ColabFold database (x-axis). In both panels, nearly all points lie above the diagonal, indicating that Logan50 yields more diverse alignments and substantially improved structural predictions.



Extended Data Fig. 6: Supporting information for identification and reactivation of HHV-6 in large-scale RNA-seq datasets. (a) Summary of HHV-6B reference transcriptome. Viral transcript abundance was computed from a prior characterization of HHV-6 reactivation in CAR T cells (Sample 34; Day 19) [22]. Boxed genes represent selected sequences used as queries for Logan-Search. (b) UMAP representation of cells profiled via scRNA-seq from lung organoid culture (PRJNA891766). Panels reflect marker genes identifying a cluster of rare proliferating T cells (1.3% of total), including two HHV-6 super-expressor cells (82% of HHV-6 UMIs). (c) Quantification of all ChIP-seq libraries from CD4+ and CD8+ TIL cultures (PRJNA901909). The abundance of HHV-6 MAPQ 30+ reads is shown with donors stratified by three participating clinical trials. Arrow indicates a high HHV-6 reactivation donor with no matched RNA-seq. (d) Single nucleotide variant analysis of Donor 10 and Donor 24 CD8+ ChIP-seq analysis. Shown are allele frequencies of 72 high-confidence single-nucleotide variants that discriminate the viral strains of the two donors.



Extended Data Fig. 7: Global distribution of AMR-associated SRA accessions (a) Summary of SRA accessions (top row) and plasmids (bottom row) categorized as AMR-positive (AMR+). First panel, amount of AMR+ vs AMR- samples in the datasets. Second panel, from the AMR+ samples, how many are classified as isolate (purple) or metagenome (yellow) as organism type. Final panel, from the AMR+ metagenome samples, distribution across metagenome categories (human: purple, soil: orange, livestock: yellow, marine: blue, freshwater: green, wastewater: red, other: grey). (b) Log2 enrichment of organism type categories in AMR+ datasets versus the average, in SRA accessions (top) and plasmids (bottom), showing relative over- or underrepresentation. Data has been randomly subsampled to avoid bias driven by categories with higher amount of data. (c) Log2 enrichment of metagenome categories among AMR+ datasets compared to the mean, for both SRA accessions (top) and plasmids (bottom). Positive values indicate overrepresentation in AMR+ samples. Data has been randomly subsampled to avoid bias driven by categories with higher amount of data. (d) Geographic distribution of unique AMR+ SRA accessions across the globe, coloured by metagenome category. Circle size indicates the number of unique accessions per location. (e) Temporal trends in AMR gene discovery. Top: Collection date timeline of AMR+ accessions by organism type (isolate: purple, metagenome: yellow). Bottom: Collection date timeline of AMR+ metagenome accessions coloured by metagenome category. (f) Distribution of AMR gene counts per accession by metagenome category. (g) Log2 enrichment of AMR gene counts per metagenome accession compared to the mean, by metagenome category. Positive values indicate metagenomes with more AMR genes per accession on average. Data has been randomly subsampled to avoid bias driven by categories with higher amount of data. In panels (c), (d), (e) bottom panel, (f), and (g), metagenome category “other” was removed from the analysis.



Extended Data Fig. 8: **Expansion of P4 phage satellite genetic diversity.** (a) Pipeline for discovering novel P4 elements. (b) Histogram of novel P4 elements binned by SatelliteFinder type. (c) Pangenome curve expressing accumulation of gene families clustered at 40% protein identity before (RefSeq: Types A, B, and C) and after Logan expansion (Logan: Types A and B). (d) Weighted Genome Relatedness Ratio Plot (wGRR) of full proteomes, defined as all proteins found between first and last detected core gene, from before (RefSeq/black) and after (Logan/grey) Logan expansion, where a darker color denotes higher similarity.