

A deep learning approach reveals unexplored landscape of viral expression in cancer

Received: 20 August 2022

Accepted: 25 January 2023

Published online: 11 February 2023

 Check for updates

Abdurrahman Elbasir¹, Ying Ye¹, Daniel E. Schäffer^{1,2}, Xue Hao¹, Jayamanna Wickramasinghe¹, Konstantinos Tsingas^{1,3}, Paul M. Lieberman¹, Qi Long³, Quaid Morris⁴, Rugang Zhang¹, Alejandro A. Schäffer⁵ & Noam Auslander¹✉

About 15% of human cancer cases are attributed to viral infections. To date, virus expression in tumor tissues has been mostly studied by aligning tumor RNA sequencing reads to databases of known viruses. To allow identification of divergent viruses and rapid characterization of the tumor virome, we develop viRNAtrap, an alignment-free pipeline to identify viral reads and assemble viral contigs. We utilize viRNAtrap, which is based on a deep learning model trained to discriminate viral RNAseq reads, to explore viral expression in cancers and apply it to 14 cancer types from The Cancer Genome Atlas (TCGA). Using viRNAtrap, we uncover expression of unexpected and divergent viruses that have not previously been implicated in cancer and disclose human endogenous viruses whose expression is associated with poor overall survival. The viRNAtrap pipeline provides a way forward to study viral infections associated with different clinical conditions.

Viral infections have a causal role in ~15% of all cancer cases worldwide¹. Viruses linked to cancer are generally divided into direct carcinogens, which drive an oncogenic transformation through viral oncogene expression, and indirect carcinogens, which may lead to cancer through mutagenesis associated with infection and inflammation. To date, seven viruses have been classified as direct carcinogenic agents in humans². Among these, the high-risk subtypes of human papillomavirus (HPV) are the causative agent of ~5% of human cancers. Chronic hepatitis B virus (HBV) or hepatitis C virus (HCV) infections are associated with most hepatocellular carcinoma cases. More recently, advances in sequencing technologies have contributed to a better appreciation of the high burden of viral infections in cancer, exemplified by Kaposi's sarcoma herpesvirus and the Merkel cell polyomavirus, which were discovered based on nucleic acid subtraction to cause Kaposi's sarcoma and Merkel cell carcinoma, respectively². The discovery of oncogenic viruses, starting with the Rous sarcoma virus³, has been critical for understanding mechanisms driving cancer

evolution and for improving cancer prevention and intervention strategies. However, the burden of viral infections in cancer is thought to remain underappreciated by much of the cancer research community⁴.

Since the advent of next-generation sequencing, new viral strains are typically identified from large-scale DNA or RNA sequencing data based on sequence similarity to known viruses. The Cancer Genome Atlas (TCGA) has become a principal resource for the identification of viral sequences in cancer tissues. Several studies screened TCGA DNA sequencing data to characterize known viruses in cancers⁵, and analyze host integration sites for viruses such as HBV that integrate into the human genome⁶. Other studies used RNA sequencing to screen for known viruses in the human transcriptome^{7–10}, and to discover novel viral isolates¹⁰. Most recently, a few studies combined DNA and RNA sequencing to quantify the presence of known cancer-associated viruses in human cancers^{11,12}. However, the set of sequenced viral clades and the set of viral clades known to infect humans are both incomplete. Viruses and cancers have rapidly evolving genomes, and a

¹The Wistar Institute, Philadelphia, PA 19104, USA. ²Computational Biology Department, Carnegie Mellon University, Pittsburgh, PA 15213, USA. ³University of Pennsylvania, Philadelphia, PA, USA. ⁴Computational and Systems Biology, Sloan Kettering Institute, New York City, NY 10065, USA. ⁵Cancer Data Science Laboratory (CDSL), National Cancer Institute, National Institutes of Health, Bethesda, MD 20892, USA. ✉e-mail: nauslander@wistar.org

new cancer-associated virus may have little sequence similarity to known viruses isolated outside of the tumor microenvironment. This issue is exacerbated when analyzing short reads, which are typical of RNA sequencing technologies. Therefore, the discovery of new and divergent cancer viruses remains highly challenging with existing strategies¹³. For the detection of bacterial viruses from metagenomic DNA sequencing, several machines and deep learning techniques have been recently developed. These methods overcome some of the limitations associated with homology-based approaches and rapidly identify viral reads including novel and divergent viruses^{14–18}. More recently, methods have been developed to identify viruses that have the potential to cause human infections^{19,20}. These recently developed methods suggest that deep learning methods to detect viral reads from RNA sequencing have the potential to uncover novel and divergent viruses in human tissues.

Here, we develop a framework, named viRNAttrap, that employs a deep learning model to accurately distinguish viral reads from RNA sequencing, and utilizes the model scores to assemble viral contigs. We apply viRNAttrap to 14 cancer types from TCGA (selected based on potential viral relevance to oncogenesis), to perform exploratory data analysis and characterize the landscape of viral infections in the human cancer transcriptome. We demonstrate the ability of viRNAttrap to identify different types of viruses that are expressed in tumors by constructing three viral databases and comparing viRNAttrap findings to sequences in those databases. We first evaluate known cancer-associated viruses that are expressed in different tumor types. Then, we curate a database of potentially functional human endogenous retroviruses (HERVs) and analyze expression patterns of different HERVs across human cancers to find that HERV expression is associated with poor survival rates. Finally, we employ viRNAttrap to identify divergent viruses that are expressed in tumor tissues. Notably, we identify a *Redondoviridae* member that is expressed in head and neck carcinomas, a *Siphoviridae* member that is expressed in 10% of high-grade serous ovarian cancers, and a *Betairedivirinae* member that is expressed in more than 25% of endometrial cancer samples. In summary, we present the first deep learning-based method to identify viruses from human RNA sequencing and demonstrate its ability to rapidly characterize viruses that are expressed in tumors and uncover viral instances that have not been previously found in these samples using alignment-based methods. viRNAttrap can be applied to identify new viruses that are expressed in a variety of other malignancies, introducing new avenues to study viral diseases.

Results

The viRNAttrap framework

To identify viruses in the human transcriptome, we first trained a neural network to distinguish viral reads based on short sequences. We collected positive (viral) and negative (human) transcripts that were segmented into 48 bp fragments and divided into training and test sets (Fig. 1a, Methods). We used different metrics to evaluate the ability of the model to identify viral sequences based on short segments. The model yielded test-set performance: area under the receiver operating characteristic curve (AUROC) of 0.81, area under the precision-recall curve (AUPRC) of 0.82 (Fig. 1b), the accuracy of 0.71, recall of 0.83, the precision of 0.67 and F1-score of 0.74 (Fig. 1c). We compared the performance of this model to previous models trained to identify viruses, namely DeepViFi¹⁶, DeepVirFinder¹⁵, ViraMiner²¹, as well as a method called “off-the-shelf Seq2Seq” compared through DeepViFi¹⁶, that does not use much domain-specific knowledge about viruses (Methods). Importantly, our model outperformed other methods in all measures, except for precision, for which DeepVirFinder outperformed all other methods (Fig. 1b, c). However, precision is less critical for this framework because alignment steps are used to further filter out negatives. Importantly, DeepViFi¹⁶, DeepVirFinder¹⁵, and ViraMiner²¹ were previously not trained or evaluated for RNA

sequencing or 48 bp reads, which is likely the reason that these methods are less appropriate in that context without specific optimization (see Methods). Examining the average model performance across segments from different human viruses, we find that human single-stranded DNA viruses from taxon *Monodnaviria* were assigned with high confidence, whereas, for RNA viruses, we observed more variation in model confidence. For example, the model confidently predicted the viral origin of sequences from Ebola and influenza viruses but assigned borderline scores to sequences from several *Phenuiviridae* members such as *Dabie bandavirus* (Fig. 1d and Supplementary Data 1).

Based on the trained neural network, we built a computational framework (Fig. 1a, Methods) to identify viral contigs from tumor RNAseq and applied the framework to 7272 samples from 14 cancer types in The Cancer Genome Atlas (TCGA)²², from which 6717 were tumor samples and 555 were non-cancer samples matched to a cancer sample from the same individual (Supplementary Data 2). In pre-processing, we extracted reads that were not aligned to the human genome (hg19) or to the phiX phage²³ that was identified as a frequent contaminant. The computational framework, named viRNAttrap, was then applied to unaligned RNA reads (to reduce the running time of viRNAttrap), to detect viral reads and assemble predicted viral contigs. Finally, in post-processing analysis, we used blastn²⁴ to compare the assembled viral contigs to three curated viral databases. We identified viral contigs originating from reference viruses that are expected in cancer tissues, human endogenous viruses, and candidate novel or more divergent viruses, which are expressed in different cancer types

Identifying reference tumor viruses

We first characterized the presence of known cancer-associated human viruses in different tumor types. High-risk human *Alphapapillomavirus* strains (HR- α HPVs) were most frequently detected; the type observed in the majority of TCGA samples is HPV16. This is expected because HR- α HPVs, such as HPV16 and HPV18, underlie ~5% of cancer cases worldwide²⁵ while low-risk human *Alphapapillomavirus* (LR- α HPV) strains, such as HPV54 and HPV201, are mostly associated with the development of genital warts but not cancer²⁶. We found at least one HR- α HPV in 288 CESC samples (286 squamous cell carcinoma samples and two non-cancer samples). We found 61 HNSC samples, and a total of 14 samples across other cancer types, that contain a contig from at least one HR- α HPV (Fig. 2a). LR- α HPVs were identified in a small set of samples mostly from matched non-cancer tissues, including cervix and head and neck (Fig. 2a and Supplementary Data 2, 3).

Hepatitis B virus (HBV) is the second most frequently detected virus across TCGA samples. HBV infections and Hepatitis C virus (HCV) infections are two primary causes of liver cancer and may co-occur in a patient¹¹. We found HBV expression in 85 LIHC tumor samples and seven non-cancer samples, and HCV in 13 LIHC tumor samples. HBV was also found in a few tumor samples and matched non-cancer samples from other cancer types (Fig. 2a). By comparing the samples predicted as virus-positive by viRNAttrap to the samples annotated as virus-positive in the TCGA clinical annotations, we found that the true positive rates of viRNAttrap were above 95% for HR- α HPVs (in CESC and HNSC), and for HCV and HBV in LIHC, supporting that viRNAttrap correctly identifies samples expressing known cancer viruses (Supplementary Fig. 1). In addition, viRNAttrap found adeno-associated virus 2 (AAV2) in eight LIHC samples, six from tumors and two from non-cancer samples. AAV2 is a small DNA virus that has the potential to integrate into human genes and contribute to oncogenesis, although the current evidence is insufficient for AAV2 to be included in the consensus list of oncogenic viruses^{27,28}. A recent study that addressed discrepancies in AAV2 expression across TCGA samples found at least one AAV2 read in 11 LIHC samples²⁷. However, in three of these samples only one AAV2 read was found, which is difficult to detect with the

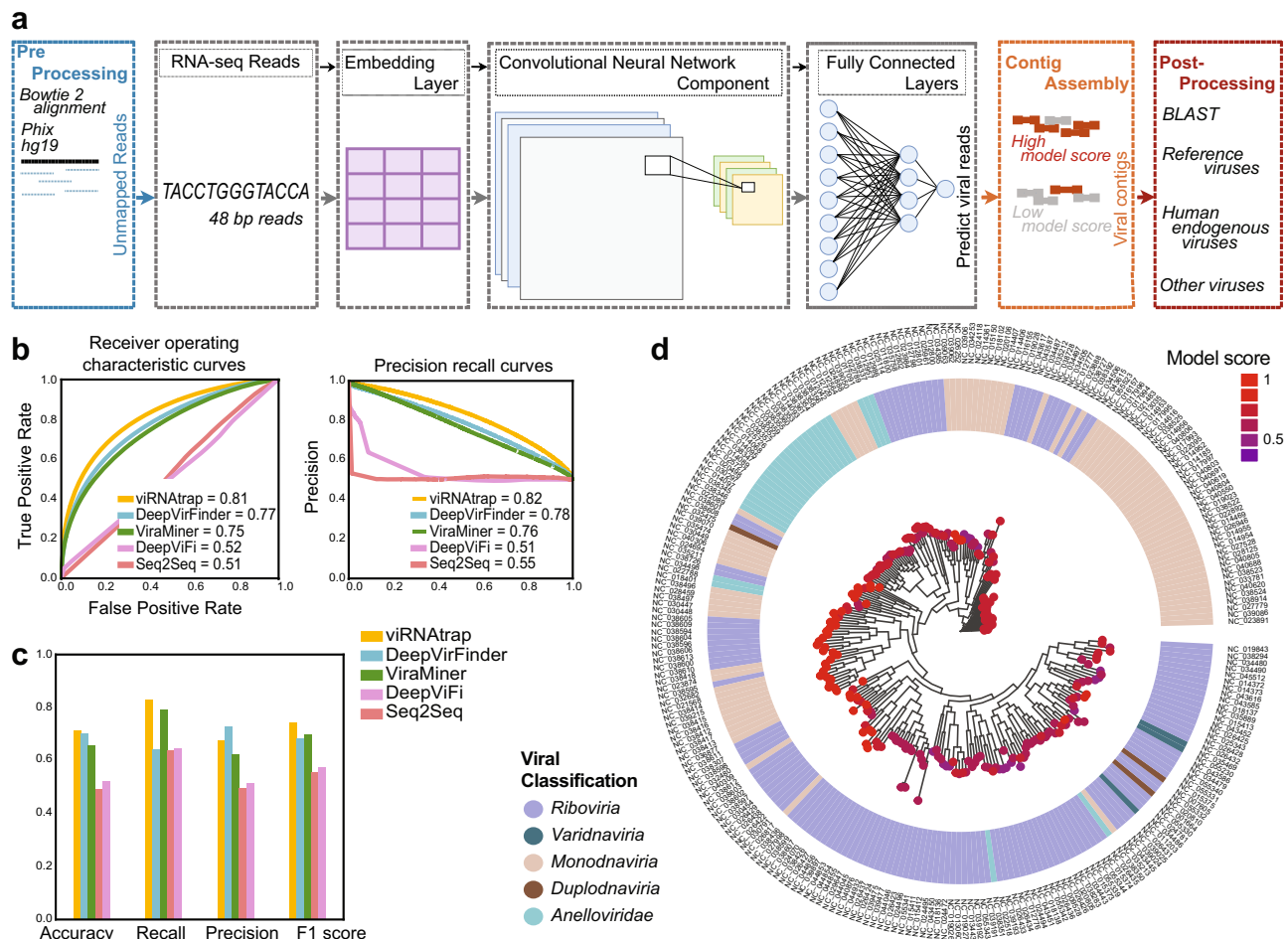


Fig. 1 | Training and evaluation of the viRNAtrap framework. **a** A schematic overview of the viRNAtrap framework. Unmapped reads were extracted and given as input to the neural network, to extract the viral reads and assemble viral contigs, that were compared against three viral databases using blastn. **b** Receiver operating characteristic and precision-recall curves showing the model performance when viRNAtrap and models used for comparison were applied to the test set. **c** Bar plots

showing different metrics to evaluate the model performance for the test set, for viRNAtrap, and models used for comparison. **d** A phylogenetic tree showing the model scores for sequences from different human viruses with the respective virus classification (using the average assigned a score for each virus). Source data are provided as a Source Data file 1.

viRNAtrap pipeline. Notably, previous studies that systematically characterized viral presence across TCGA did not identify AAV2 in more than six LIHC samples^{11,27}, demonstrating the sensitivity of viRNAtrap compared to other computational methods. We additionally detected AAV2 in one KIRC sample, one PAAD sample, and one matched non-cancer sample from LUAD (Fig. 2a).

We found several samples that express human polyomaviruses, especially polyomaviruses 6 and 7. Most notably, we found seven BRCA samples and two HNSC samples that express polyomaviruses. We additionally found Parvovirus B19 sequences in a few samples²⁹ (three cancer and one matched non-cancer); this virus has been mostly associated with normal tissues³⁰, but was also previously identified in isolated tumor cases^{31,32}. We investigated possible genomic correlates of the expression of these viruses, including the tumor mutation burden (TMB, the rate of somatic mutations in a tumor, which is a biomarker and is annotated for all TCGA samples), and the chromosome-level aneuploidy (Methods). We found that HR- α HPV-positive samples have lower TMB and aneuploidy levels compared to HR- α HPV-negative samples (Fig. 2b). In contrast, LIHC cancer patients positive for HBV showed significantly higher TMB compared to HBV-negative samples (Supplementary Fig. 2). We additionally examined the association between the expression of known oncoviruses and overall survival. While none of the associations were significant after adjustment for

multiple hypotheses (Supplementary Fig. 2 and Supplementary Table 1), we found a trend that HR- α HPV-positive HNSC patients have better survival compared to HR- α HPV-negative patients (by the Kaplan–Meier curves Fig. 2c), which is confirmatory of previous studies^{33,34}. We also found a positive association between the viral presence and the overall survival of LIHC patients with HBV (Supplementary Fig. 2 and Supplementary Table 1).

Uncovering expression patterns of HERVs in cancer tissues

To further demonstrate the utility of viRNAtrap, we analyzed the expression of HERVs across different tumor types in TCGA (HERVs were not used to train the viRNAtrap model). HERVs constitute ~8% of the human genome; most HERV sequences are remnants of ancestral retroviral infection that became fixed in the germline DNA^{35,36}. HERV proteins are found expressed in different conditions including cancer tissues^{37–41}. Specifically, the HERV-K family, which was most recently integrated into the human genome and is one of the most abundant HERV families in the human genome (along with HERV-H), was previously reported in tumor tissues and cell lines^{42,43}. Moreover, recent findings reported the association between HERV expression and poor survival rates^{42,36,44–46}.

To comprehensively characterize HERV members that are expressed in different tumors, we established a database of potentially

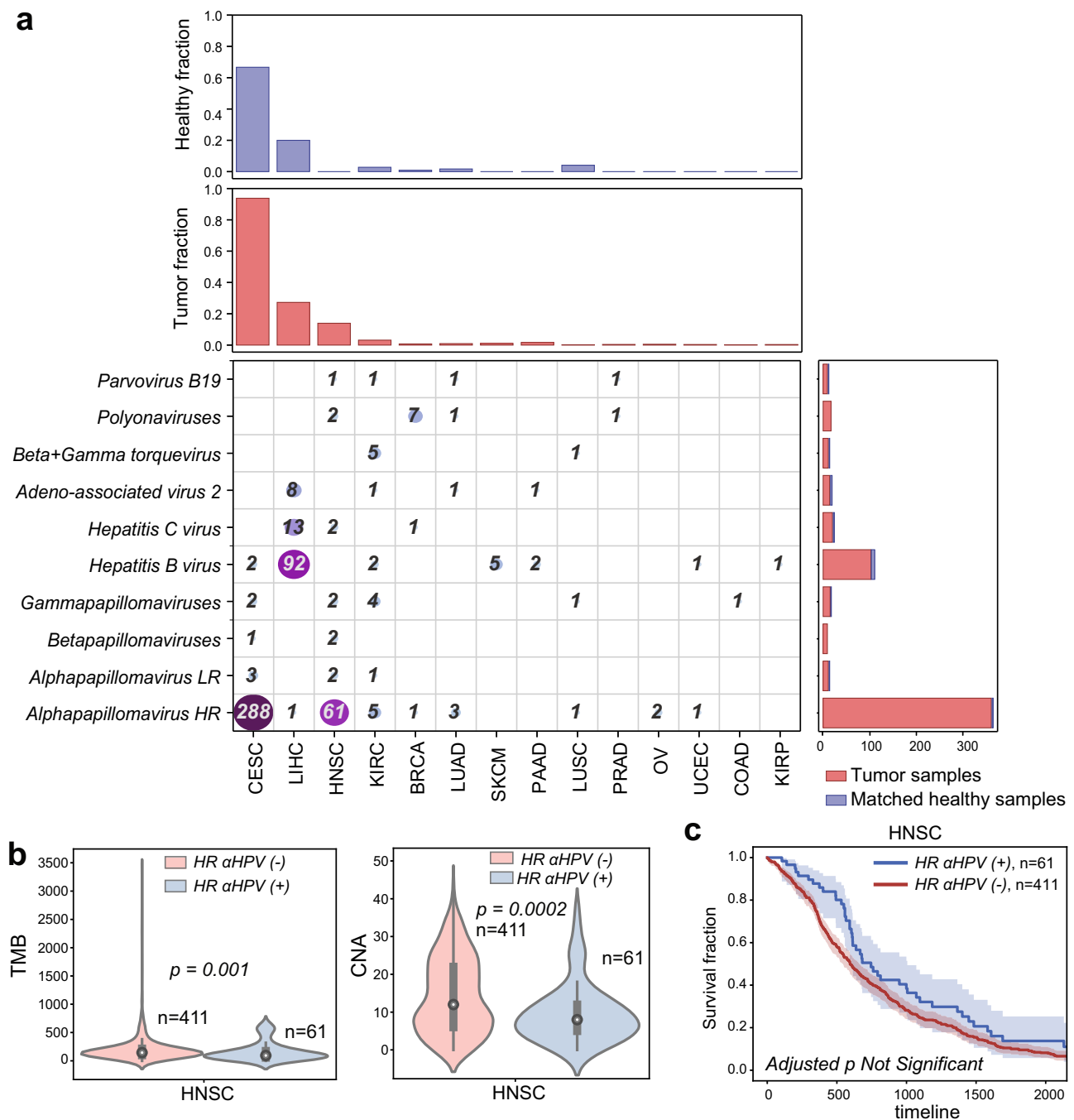


Fig. 2 | Reference human viruses expressed in different tumor types. a Heatmap showing the total number of virus-positive samples identified from RNA sequencing in different tumor tissues. The top panels show the fraction of tumor and non-cancer samples in which viruses were identified. The right panels show the number of viruses found in tumor and non-cancer samples. **b** Violin plots comparing the tumor mutation burden (TMB) and the number of chromosome-level copy number alteration (CNA) between HNSC patients where expression of high-risk alpha papillomaviruses was detected vs those patients where expression of high-risk alpha papillomaviruses was not detected. Black dots represent the medians, and

the boundaries of the violin plots refer to the maximum and minimum values, respectively. Two-sided Wilcoxon rank-sum *p* value is reported. **c** Kaplan–Meier curves comparing the survival rates between HNSC patients where the expression of high-risk alpha papillomaviruses was detected (blue curve) vs those where the expression of high-risk alpha papillomaviruses was not detected (red curve). The FDR-adjusted two-sided log-rank *p* value is not significant (Supplementary Table 1). For Kaplan–Meier curves, shaded areas represent the confidence interval of survival. Source data are provided as a Source Data file 2.

functional HERVs that were extracted from the human genome (Methods). The viRNAtap contigs were aligned against this database, to identify patterns of HERV expression in the 14 cancer types considered throughout this study.

As expected, we found that the most abundantly expressed HERV families are HERV-K and HERV-H. The fraction of samples expressing

different individual HERV members was used to cluster tumor types. Interestingly, we found that squamous cell carcinomas (including cervical, lung, and head and neck) are clustered together based on the proportional distribution of expressed HERV members (Fig. 3a). The HERVs that are most abundantly expressed across different cancers include some that are in proximity to cancer-associated genes or single

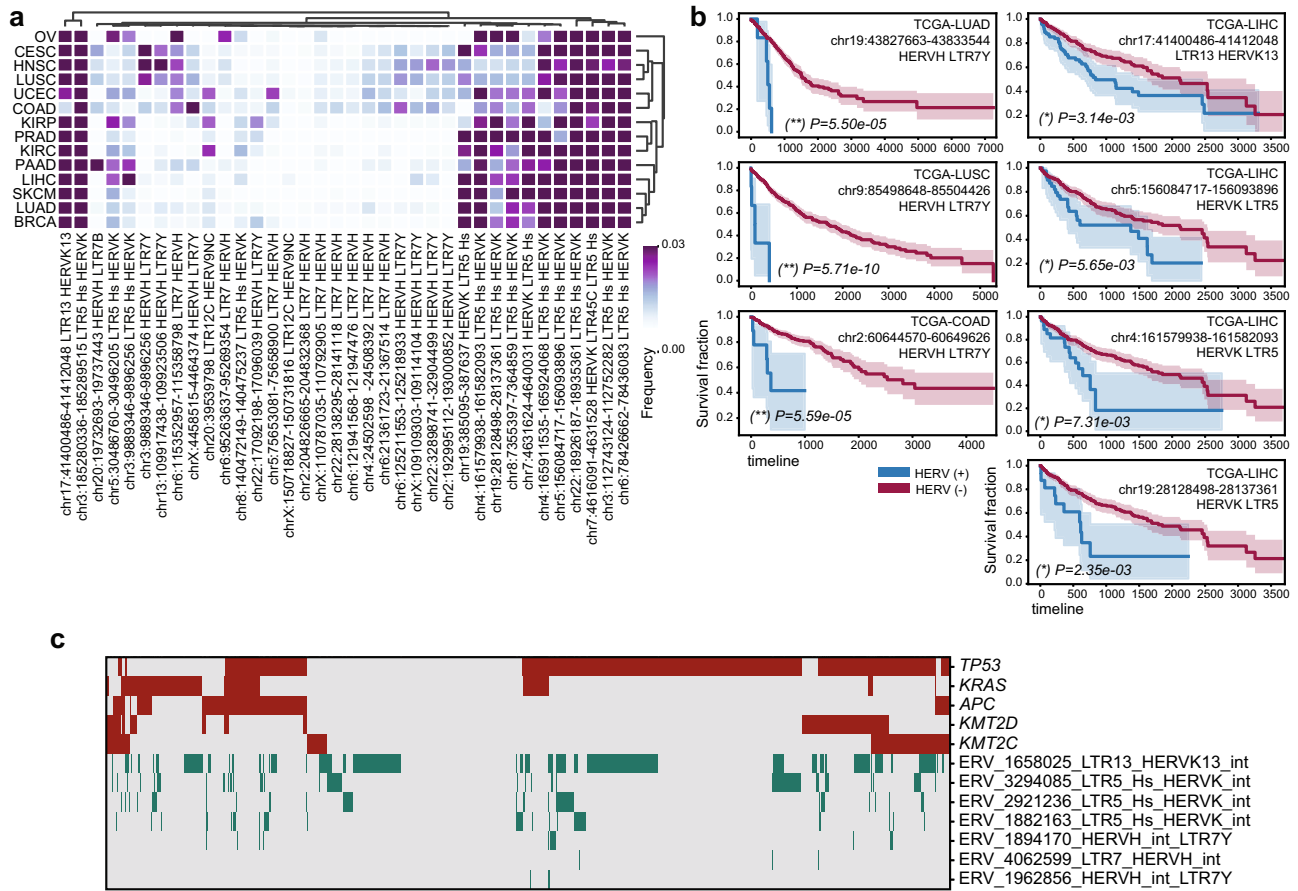


Fig. 3 | Human endogenous retroviruses (HERVs) expressed in different cancer types. **a** Heatmap clustering the proportion of HERVs across different tumor types. The rows are 14 TCGA tumor types. The 36 columns are the 36 distinct HERVs with the highest expression in human cancers, mapped to unique regions in the genome (Supplementary Data 5). **b** Selection of Kaplan–Meier curves comparing the survival rates between patients in which any HERV reads were detected (blue curves) versus those in which no HERV reads were detected (red curves). The unadjusted two-sided log-rank *p* values are reported. (**) global FDR *q* < 0.05, (*)

cancer-type specific FDR *q* < 0.05. For Kaplan–Meier curves, shaded areas represent the confidence interval of survival. Additional significant associations between HERV and survival are reported in Supplementary Data 12. **c** Heatmap showing somatic mutations in major cancer driver genes (selected are the most frequently mutated driver genes in these samples, red) and the expression of HERVs that are significantly associated with survival in LIHC, LUAD, LUSC, and COAD (green). Source data are provided as a Source Data file 3.

nucleotide polymorphisms (SNPs) (Supplementary Data 3, 4). Specifically, one HERV-H member (chr2:204826665-204832368) is located 365 bp from the *ICOS* (Inducible T-cell costimulatory) gene, which has been associated with tumor immune responses^{47–50}. In addition, one HERV9 member (chrX:150718827-150731816) is located 330 bp from the *PASDI* cancer/testis antigen gene (each of these two HERVs are found in ten TCGA samples, Supplementary Data 4, 5).

We investigated associations between HERV transcript presence and patients’ overall survival (Fig. 3b). In agreement with previous studies^{12,36,44–46}, we find that patients with HERV-K- and HERV-H-positive cancer samples have significantly lower overall survival compared to HERV-K- and HERV-H-negative patients in COAD, LUSC, LUAD, and LIHC. Notably, every significant association that we identified between HERV presence and overall survival in these cancer types is negative (Fig. 3b and Supplementary Table 2).

To investigate the link between HERV expression and poor survival, we compared the TMB and aneuploidy scores between patients expressing HERVs and those without HERV expression. HERVs that were associated with poor survival were not associated with TMB or aneuploidy (Supplementary Data 6). We found that HERVs associated with poor overall survival were generally more likely to be expressed in the presence of somatic mutations in frequently mutated cancer driver genes, such as *TP53*, *KRAS*, *ARID1A*, and *PTEN* (using hypergeometric

enrichment, Supplementary Data 7). However, we did not find a strong association with mutations in any specific gene, and HERV expression was found even in samples with no somatic mutations in any of these genes (Fig. 3c and Supplementary Data 8)

Finding divergent viruses in human cancer

We next investigated tumor expression of divergent viruses that have rarely or never been previously reported in human cancers. We aligned the contigs produced by viRNAtrap against a database of viruses (Methods) from different hosts that were not expected to be found in tumor tissues, including human, bat, mouse, insect, plant, and bacterial viruses. (Fig. 4a). We found multiple contigs of mosaic plant viruses in distinct samples from most tumor types, especially adenocarcinomas. For example, the watermelon mosaic virus was found in three colorectal cancer samples, and the Bermuda grass latent virus, which was previously reported in a COAD sample¹⁰, was identified in multiple samples from three cancer types (COAD, LIHC, and UCEC; Fig. 4a). Mosaic plant viruses have been previously detected in human feces^{51,52}, which could suggest viral entry and travel through the digestive tract. However, it is unclear how mosaic plant viruses would reach other tumor tissues, such as the liver and the endometrium, and whether these are associated with an unidentified source of laboratory contamination.

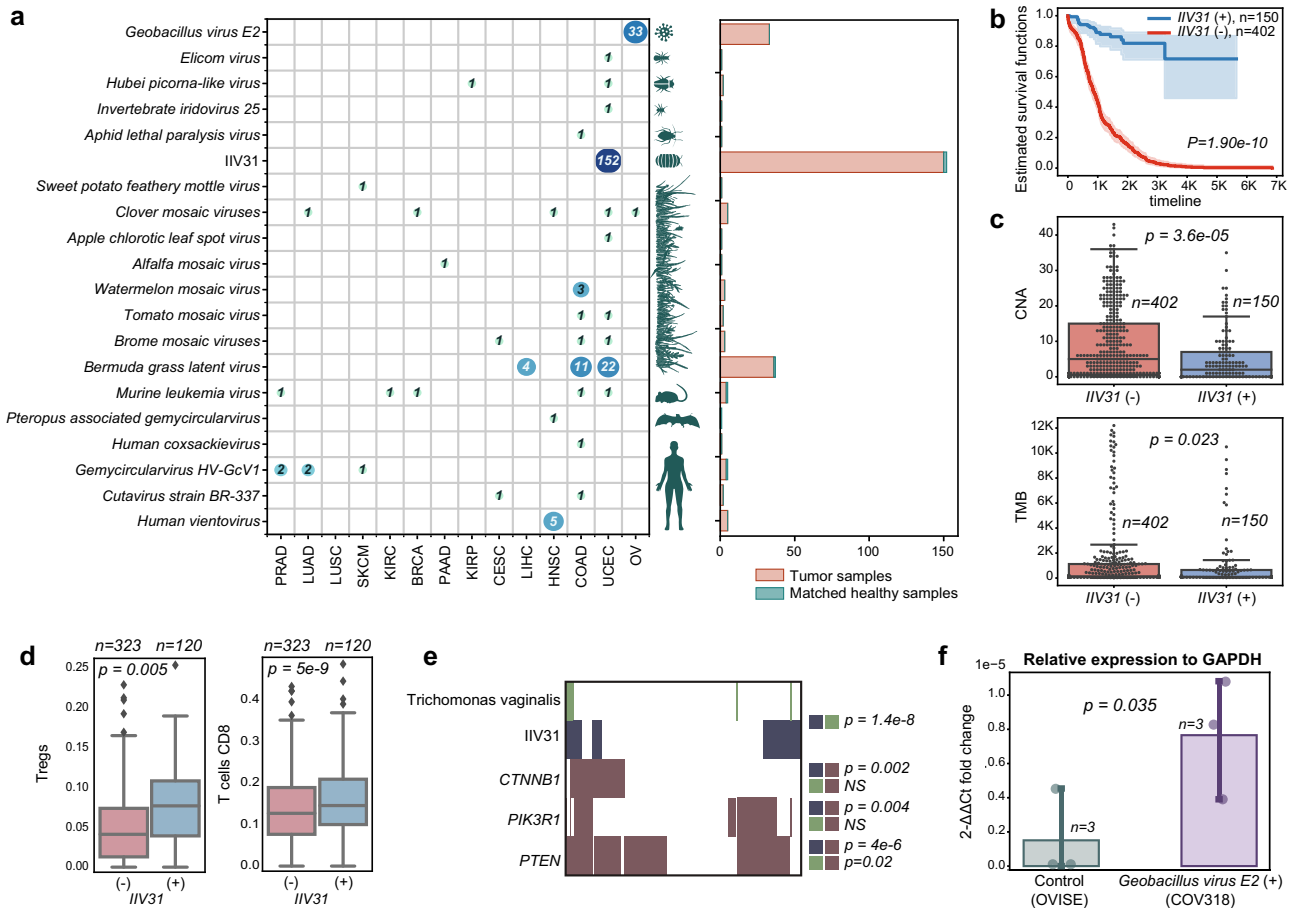


Fig. 4 | Unexpected and divergent viruses infecting different host taxa across TCGA samples. **a** Unexpected and divergent viruses expressed in TCGA samples. Each row in the matrix represents one virus and the entry in each column indicates the number of cancer samples of each type in which each virus was detected. The canonical hosts of each virus are depicted at the left of the matrix. At right, the aggregate number of tumor and normal samples containing reads of each virus are shown in a bar plot. **b** Kaplan–Meier curves comparing the survival rates between patients in which IIV31 reads were detected (blue curves) vs those where viral reads were not detected (red curves). For Kaplan–Meier curves, shaded areas represent the confidence interval of survival. The log-rank p value is reported. **c** Box plots comparing the chromosome-level copy number alteration (CNA, top panel) and the tumor mutation burden (TMB, bottom panel) between cancer patients where IIV31

is found (blue) and patients where IIV31 is not found (red). Two-sided Wilcoxon rank-sum p value is reported. **d** Box plots comparing CIBERSORT-inferred proportions of regulatory T cells (Tregs) and CD8 T cells between patients positive and negative for IIV31. Boxes show the quartiles (0.25 and 0.75) of the data, center lines show the medians, and whiskers show the rest of the distribution except for outliers. A two-sided Wilcoxon rank-sum p value is reported for comparisons assigned with $FDR < 0.05$. **e** *Trichomonas vaginalis* and somatic mutations in *PTEN*, *CTNNB1*, and *PIK3R1* are associated with IIV31 presence. One-sided Fisher’s exact test p values are provided. **f** Bar plot comparing the fold change (relative to GAPDH) between the COV318 cell line that was predicted as *Geobacillus*-positive, and the OVISE cell line that was used as control. Error bars show the standard deviation. The one-sided t -test p value is provided. Source data are provided as a Source Data file 4.

Notably, we identified expression in five head and neck carcinoma samples of a *Vientovirus*, a member of the recently characterized human virus family *Redondoviridae* that is associated with the human oro-respiratory tract⁵³ (Fig. 4a and Supplementary Data 3, 9). We also found expression of a *Gemycircularvirus* HV-GcV1⁵⁴ in distinct samples from several cancer types, and *Cutavirus* expression in one COAD and one CESC sample each. We additionally detected human coxsackievirus⁵⁵ in a COAD sample, confirming a previous report¹⁰.

We also found expression of a few arthropod viruses in TCGA, almost exclusively in UCEC samples (Fig. 4a), the most notable of which is *Armadillidium vulgare* iridescent virus (IIV31)⁵⁶. We detected reads that align with IIV31 proteins in 152 endometrial cancer samples (which constitute more than 25% of endometrial cancer samples studied). While we did not find previous reports of IIV31 in these samples, reads that align to the same strain were recently detected in a few DNA sequencing samples, but were filtered because these were not included in databases of multiple pipelines¹². IIV31 is in *Betairidovirinae*; members of this subfamily of dsDNA viruses infect a wide variety of arthropods, including common insect parasites of humans⁵⁷. One

study speculated on the role of *Betairidovirinae* transmitted by mosquitoes in human disease⁵⁸, but, to our knowledge, their presence in humans has not been reported before. While *Betairidovirinae* are not considered to be pathogens of vertebrates, one study showed that the model *Betairidovirinae* insect iridovirus 6 (IIV6) was lethal to mice after injection, while heat-inactivated IIV6 was not⁵⁹. Additional studies have shown that *Betairidovirinae* can infect vertebrate predators of infected insects as well as several vertebrate cell lines⁶⁰. Therefore, *Betairidovirinae* may opportunistically infect vertebrates, including humans.

We identified different IIV31 genes expressed in UCEC samples, and samples positive for IIV31 proteins originate from different batches and sequencing centers (Supplementary Data 10). In addition, we found that IIV31 presence was strongly and positively associated with overall survival (Fig. 4b), and negatively associated with TMB and chromosome-level aneuploidy (Fig. 4c, d). We did not identify a path to contamination by IIV31; the multiple origins of IIV31-positive samples and significant associations between IIV31 expression and other cancer properties both suggest that IIV31 is not a contaminant. Of the most highly expressed IIV31 proteins, we found an IAP apoptosis inhibitor

homolog and serine/threonine protein kinases that were individually associated with poor overall survival (YP_009046765, YP_009046752, and YP_009046774, respectively), as well as a *RAD50* homolog (YP_009046808, Supplementary Fig. 3 and Supplementary Data 10).

We found a significant positive association between IIV31 and CIBERSORT⁶¹ inferred CD8⁺ T-cell frequency and Treg frequency (Fig. 4d). These findings, together with the association with improved survival, suggest that IIV31 could be linked with a different infection, either directly or indirectly. We explored the association of IIV31 infection with *Trichomonas vaginalis* (TV)⁶² infection. TV is a single-celled protozoan pathogen that infects the human urogenital tract⁶³, and has been associated with an increased risk of cervical cancer, which is enhanced by HPV coinfection⁶⁴. We found that TV is expressed in multiple UCEC tumor samples (we verified 21 TV-positive tumors with strict alignment parameters, due to a high false positive rate when aligning against TV transcripts). Indeed, TV-positive samples are highly enriched with IIV31-positive samples (Fisher exact test p value = 1.4e-8). Both TV and IIV31 are significantly associated with somatic *PTEN* mutations, which are linked to better survival in endometrial cancers⁶⁵ (whereas the presence of IIV31 is also associated with mutations in *CTNNB1* and *PIK3RI*, Fig. 4e).

We additionally identified *Geobacillus* virus E2 expression in 33 ovarian cancer samples; this virus is likely the most frequently expressed virus in high-grade serous ovarian cancer. To further validate the presence of the *Geobacillus* virus E2, we applied viRNAtrap to cell line data from CCLE⁶⁶. We identified the COV318 cell line as *Geobacillus* virus E2-positive and identified the OVISe cell line as a virus-negative control. Through qRT-PCR we validated the expression E2 in the predicted-positive cell line COV318 (Fig. 4f). These results verify that *Geobacillus* virus E2, which was never found in ovarian cancer before, is indeed expressed in ovarian cancer cells, and that viRNAtrap can be used to sensitively detect virus-positive samples. *Geobacillus* bacteria has been previously detected in multiple ovarian cancer samples^{67,68}. While we could not pinpoint the *Geobacillus* species harboring the phage, it is likely within those previously found in ovarian cancer samples^{67,68}.

We found murine leukemia virus⁶⁹ expression in distinct samples from five cancer types. However, murine leukemia virus contamination has been reported for cell culture due to human DNA preparation⁷⁰. Our method additionally detected a previously unknown virus in a matched non-cancer sample from one HNSC patient, with protein similarity to *Pteropus* (fruit bat)-associated *Gemycircularvirus* and several other gemycircularviruses (Supplementary Data 3, 9).

Discussion

Identification of viruses from tumor RNA sequencing allows for the potential discovery of new carcinogenic agents and mechanisms. The discovery of novel and divergent viral species that contribute to cancer initiation and progression is crucial for the development of new therapeutics, including vaccinations, screening practices, and antimicrobial treatments. Viruses are currently identified from sequencing reads based on similarity to known viruses⁷¹. However, when studying viruses from short reads, typical with Illumina-based RNA sequencing, reads originating from divergent viruses may share little sequence similarity to known viruses, rendering the identification of novel viruses highly challenging.

To address this challenge, we developed viRNAtrap, a new, alignment-free framework to identify viral reads from RNAseq and assemble viral contigs. The contigs detected by viRNAtrap can be aligned to different viral databases, as we demonstrate in this study, to rapidly identify viral expressions of interest in tumor samples. We curate a database of HERVs that comprise intact retroviral genes in the human genome and survey the expression of these viruses across different cancer tissues. Through a database of divergent viruses, we demonstrate that viRNAtrap identifies viruses in TCGA samples that

were not detected in previous studies. This is enabled through an integrative method that uses the model scores to assemble viral reads rather than aligning short divergent reads to viral databases or applying assembly to many unmapped reads. We further show that using the deep learning model substantially improves the running time, while not compromising sensitivity if more than five viral reads are present (Supplementary Fig. 4, see Methods). Importantly, the output of viRNAtrap can be alternatively used as input to motif search tools, to potentially identify highly divergent viruses. Because the deep learning model underlying viRNAtrap was trained to distinguish viral from human sequences, the model predictions for sequences derived from a range of other organisms is not defined. Future work could train models to identify viruses from a variety of other organisms, and, with the viRNAtrap framework, achieve higher sensitivity for viral detection.

We employ viRNAtrap for exploratory data analysis and characterize viruses that are expressed across 14 cancer tissues from TCGA and analyze their genomic and survival correlates. Interestingly, while the expression of some exogenous cancer viruses is known to be associated with improved survival, we found that the expression of human endogenous viruses is strictly associated with poor survival rates. Expression of a virus of the subfamily *Betaireidovirinae*, which are pathogens of insects, found in endometrial cancer tissues was similarly associated with significantly better overall patient survival. For all divergent viruses reported in this study, the presence and classification of multiple viral reads was verified by targeted blastn- and blastx-based sequence analyses in different samples. However, it is not possible to model all contaminants of viruses that may have infected the samples during laboratory procedures¹⁶.

Perhaps, the most interesting divergent virus we found is IIV31 from the subfamily *Betaireidovirinae*, which was frequently detected in UCEC TCGA samples. Interestingly, IIV6, a very close relative of IIV31, can infect a variety of vertebrates including mice, and induces an immune response in mammalian tissues^{60,72}. Thus, one possibility is that IIV31 is transmitted to the uterus through another insect, such as the crab louse. While we have not yet confirmed the source of this virus, our results imply that its presence may be a direct or indirect consequence of *Trichomonas vaginalis* infection. Therefore, it shows that viRNAtrap is sufficiently powerful to identify a previously unknown viral transcript in tumor samples, whether oncogenic or neutral. Through this analysis, we also identified TV reads in multiple endometrial cancer samples, indicating a possible new association between TV and endometrial cancer, like the known association of TV with cervical cancer⁶⁴. One of the established pathogenic mechanisms of TV infection in humans, which may also explain the frequent HPV coinfection, is that TV secretes exosomes that have the effect of suppressing CXCL8⁷³. Interestingly, low expression of CXCL8, like infection with TV, has been associated with a favorable prognosis in cervical cancer⁷⁴. Thus, it is possible that the presence of IIV31 is a secondary infection in patients already infected with TV or some other pathogen that suppresses the human anti-viral response.

Importantly, we identified the E2 *Geobacillus* virus in 10% of high-grade, serous ovarian cancers, making it the most frequently expressed virus in this cancer type. We experimentally verified that E2 *Geobacillus* is indeed expressed in cell lines. We also found expression of a *Redondoviridae* member in head and neck cancers that was not previously reported⁷⁵. This finding calls for a study of the role of *Redondoviridae* in tumor initiation and progression, as this family of viruses was only recently detected in humans and associated with different clinical conditions.

In conclusion, we developed viRNAtrap, a new software for alignment-free identification of viruses from RNAseq, allowing rapid characterization of viral expression and detection of divergent viruses. We applied it to tumor tissues from TCGA, uncovering expression patterns of different groups of viruses. We report previously unrecognized associations between several forms of cancer and several

unexpected viral clades, including viral clades canonically found in produce and in insect parasites of humans. Future studies may employ viRNAtrap to find viruses that contribute to other malignancies.

Methods

Training a neural network to distinguish viral RNA sequencing reads

The viRNAtrap framework is composed of two main components, illustrated in Fig. 1a. The first is a deep learning model, which was trained to accurately distinguish viral from human reads using RNA sequencing. The second assembles the predicted viral reads into contigs. The trained neural network is composed of one 1D-convolutional layer and three fully connected layers, one of which is the final output layer. The RNA sequences were one-hot encoded to vectors that were given as input to the model. The learning rate was set to 0.0005, we used 64 filters with ReLU as an activation function in the convolutional layer, followed by one pooling layer for feature extraction. The global extracted features from the convolutional layer are passed to three fully connected layers, to make a prediction based on a sigmoid activation function in the output layer.

To train the model, we collected human and viral sequencing data. Coding sequences of human and other placental viruses were downloaded from the Virus Variation Resource⁷⁶. Human transcripts for hg19 were downloaded from NCBI Human Genome Resources⁷⁷. These sequences were segmented into 48 bp segments, which is the read length for the RNAseq in almost all tumor types in TCGA; only a few tumor types that were added chronologically last to TCGA used longer reads. We used a 48 bp window size for human transcripts and a 2 bp window size for viral sequences, to balance the positive and negative data. Then, these were randomly split (where all segments of each transcript were considered together) into balanced train, validation, and test sets ($n = 8,000,000, 800,000, \text{ and } 2,558,044$, respectively).

Model performance evaluation and comparison to existing methods

We evaluated the performance of the model using the area under the receiver operating characteristic curve (AUROC), the area under the precision-recall curve (AUPRC), as well as accuracy, precision, recall, and F1-score, for the test dataset. We trained multiple models with different architectures and hyperparameters and then selected the model with the highest average between the validation-set AUROC and recall. The model was trained using TensorFlow 2.6.0 and Keras⁷⁸. We compared the performance of our model to models from DeepViFi¹⁶, DeepVirFinder¹⁵, ViraMiner²¹, and off-the-shelf Seq2Seq model. Because this is the first approach trained to predict viruses from RNA sequencing reads of length 48 bp, we used our training data to retrain each of these models, following the instructions provided by each method, and evaluated the AUROC, AUPRC, accuracy, precision, recall, and F1-score using our test set (see Supplementary Methods for a detailed description of hyperparameters used). Importantly, existing methods were not designed for reads shorter than 150 bp, therefore they should not be expected to perform as well as viRNAtrap on 48 bp segments, for which viRNAtrap was optimized. Our comparison does not rule out the possibility that new hyperparameter optimization for this purpose may enhance the performance of existing methods for 48 bp sequences.

Assembling viral contigs from neural network predicted viral reads

Once the viRNAtrap model predicts the probability of a viral origin of each read, reads with model scores more than 0.7 are used as seeds to assemble viral contigs. Viral contigs are assembled using an iterative search for substrings with exact matches between 24 bp k-mers. Each seed is complemented from the left and right ends using its left-most and right-most 24 bp k-mers. For both the left and right assembly,

reads containing the left or right-most k-mers in a different position from the read that is being searched are identified. The read adding the maximal number of bases to the assembled contig is used to complement the left and right contigs. The model scores that were assigned to reads that are used to assemble each contig were averaged, and the assembly terminates if the average score is below 0.5. Finally, the right and left contigs are concatenated, to yield a complete viral contig. This algorithm was implemented in Python 3 and subsequently in C, which improved the running time by more than an order of magnitude for inputs with large numbers of reads.

Data pre-processing

We downloaded RNA sequencing data from Genomic Data Commons (GDC; <https://portal.gdc.cancer.gov/>)⁷⁹ as BAM files. High-quality reads were selected and mapped with Bowtie2 against hg19 (1000 Genomes version) and PhiX phage (NC_001422), and only the unmapped reads were kept. Then, we merged the paired-end reads and converted them to fastq files, which were used as input for the viRNAtrap framework, to yield predicted viral contigs.

Viral databases

Viral contigs yielded by the assembly component were used as inputs to blastn²⁴. Three databases were used to search for viruses (with an E-value threshold of 0.01):

(1) RefSeq reference human viruses, downloaded from the National Center for Biotechnology Information (NCBI)⁷⁷, to which we added human papillomaviruses strains that are not in RefSeq from PAVE (<https://pave.niaid.nih.gov/>)⁸⁰. Reference viruses were searched using blastn, with default parameters except for a word size of 15 (lower than the default of 28), which was chosen to allow identification from short contigs.

(2) more divergent viruses were obtained from RVDB⁸¹ (<https://hive.biochemistry.gwu.edu/rvdb/>) which was then filtered to remove non-viral elements, endogenous viruses, and accessions that were consistently not verified using blastn against the nonredundant (nr) blast nucleotide database.

(3) Human endogenous viruses. We curated a database of potentially functional HERVs through the evaluation of viral protein completeness (in contrast to a previous study that evaluated HERV expression in distinct RNAseq datasets⁸²). The initial genomic locations of reported HERV elements were downloaded from the HERVd HERV annotation database (<https://herv.img.cas.cz/>)⁸³. The nucleotide sequences in hg19 for each reported HERV were extracted using twoBitToFa⁸⁴. We then applied blastx against NR with an E-value cutoff of 1E-4, as well as a profile search⁸⁵ against collected POL proteins, where the profile was obtained by collecting POL genes annotated in GenBank in lentiviruses (as of September 2016) and aligning their amino acid sequences using MAFFT⁸⁶. Sequences with at least one identified retroviral protein motif of POL/RT, GAG, or ENV were extracted, yielding 3044 HERVs that were considered for search in TCGA samples (Supplementary Data 5). Importantly, the high mutation rate of HERV⁸⁷ prohibits most HERV sequences from aligning to the human genome in pre-processing^{12,88}, however, in rare cases, HERV regions that are conserved would not be identified by this approach.

Quality standards for virus identification

For all viruses, blastn was applied with an E-value cutoff of 0.01 and any sequences with a match to contaminant accessions (that were associated with vector contamination) were filtered out.

a. Reference viruses. For every sample, contigs mapped to each accession were extracted. Identified accessions with maximum qcov across contigs of more than 90%, average qcov of more than 50%, and average similarity of more than 90% were considered. Accessions with maximal contig length under 100 bp were manually inspected and verified against nr.

- b. Human endogenous viruses. For every sample, contigs mapped to each HERV were extracted. HERVs with contigs longer than 200 bp, and with average qcov and similarity of more than 95% were considered.
- c. Divergent viruses. For every sample, contigs mapped to each accession were extracted. Viruses already identified through the reference database were removed. Identified accessions with maximal contig length of more than 300 bp and qcov of more than 40%, or with maximal contig length of more than 100 bp and qcov of more than 75% and average similarity of more than 75% were considered for manual inspection.

All instances of divergent viruses identified in TCGA samples were verified using `blastn` against `nr`, to support that the virus strain is indeed the best match to a viral contig generated by `viRNAtrap`. We reason that non-reference viruses (divergent viruses and viruses of non-human hosts) that were identified and verified in more than one sample were less likely to be a contaminant or isolated events, whereas samples with fewer reads from such viruses may be filtered due to the strict filtering. We therefore additionally searched using the STAR aligner³⁹ across tumor types where these viruses were identified through `viRNAtrap` (Supplementary Data 3). The following accessions were additionally searched using STAR to increase sample coverage (as these were the most interesting divergent strains found across multiple samples): Bermuda grass latent virus (NC_032405), *Armadillidium vulgare* iridescent virus IIV31 (NC_024451), *Geobacillus* virus (NC_009552), and the human lung-associated virovirus (NC_055523).

Filtering contaminants

To filter vector contaminants, we applied `VecScreen`⁹⁰ to the assembled contigs that have been mapped to viruses through our databases, where virus accessions associated with vector contaminants were entirely removed from the search (Supplementary Data 11).

In addition, we examined the application of software such as `Kraken2`⁹¹ to the RNAseq reads for filtering reads that are not likely of viral origin, by applying `Kraken2` to reads of LIHC samples. However, we found that 99% of the reads would not be filtered using this approach (Supplementary Fig. 5), likely due to the short reads (48 bp) for which `Kraken` has not been designed or evaluated, as longer sequences are known to be more accurately mapped⁹².

Genomic correlates of viral expression

We correlated viral expression with genomic markers across TCGA samples. Chromosomal aneuploidy levels for TCGA samples were extracted from⁹³ and the total number of chromosome-arm-level alterations was used. The tumor mutation burden was defined to be the total number of somatic mutations in each sample, downloaded from the Xena browser⁹⁴ (<https://xenabrowser.net>). `CIBERSORT`⁶¹ software was applied to TCGA samples using the default set of 22 immune-cell signatures.

Cells and culture conditions

Human ovarian cancer cell lines COV318 and OVISE were cultured in RPMI1640 medium containing 10% fetal bovine serum (FBS) and 1% penicillin-streptomycin under 5% CO₂. All of the cell lines were authenticated at The Wistar Institute's Genomics Facility using short-tandem-repeat DNA profiling. Regular mycoplasma testing was performed using a LookOut mycoplasma PCR detection kit (Sigma, cat. no. MP0035).

Experimental validation of the *Geobacillus* virus E2 in ovarian cancer cell lines

Reverse-transcriptase qPCR (RT-qPCR) RNA was extracted using TRIzol reagent (Invitrogen, cat. no. 15596026). Extracted RNA was used for

reverse-transcriptase PCR using a High-capacity cDNA reverse transcription kit (Thermo Fisher, cat. no. 4368814). Quantitative PCR was performed using a QuantStudio 3 real-time PCR system. GAPDH was used as an internal control. The fold change was calculated using the 2- $\Delta\Delta C_t$ method. The primers used for reverse-transcriptase qPCR are: GAPDH forward, GTCTCTCTGACTTCAACAGCG and reverse, ACCACCCTGTTGCTGTAGTAGCCAA; *Geobacillus* virus E2 terminase forward, TTGCGATGCGTACTCAGACT and reverse, CTCTTTTTGGTCAGCAGCGG. Primers were obtained using NCBI primer design tool as shown in the Supplementary Information. The primers were synthesized by Integrated DNA Technologies IDT. A specification of the primer design is provided in the Supplementary Information.

Identification of *Trichomonas vaginalis*-positive samples

UCEC unmapped (to hg19) reads were aligned to the reference genome of *Trichomonas vaginalis* (GCF_000002825)⁶² strain G3 using `blastn`²⁴ with E-value <1e-8 and more than 90% identity. These thresholds were set to remove false positives that were frequent when aligning against *Trichomonas vaginalis* when examining both `blastn`²⁴ and STAR aligner³⁹. TV reads for each TV-positive sample were verified by manual inspection of the output alignments.

viRNAtrap performance evaluation

To evaluate the contribution of the model to the `viRNAtrap` pipeline we re-ran `viRNAtrap` on 10 LIHC samples, and additionally ran a modified `viRNAtrap` pipeline not using the model, on the same system. We compared the viruses identified by the model-based approach to those that are identified when the pipeline is applied without using the model (Supplementary Table 2, showing similar viruses with a different number of contigs). We additionally compared the running time of the two approaches (Supplementary Fig. 4).

To evaluate the sensitivity of the `viRNAtrap` pipeline based on the number of viral reads present in a sample, we performed a simulated analysis. From the test dataset, we downsampled groups of viral reads with different group sizes (10,000 groups for each size, from one read up to 10 reads), and we evaluated the number of groups with at least one read that is scored above 0.7, which is the seed threshold used for the `viRNAtrap` assembly. Therefore, this analysis is estimating the probability of identifying viruses based on the number of reads present. We found 93 and 99% of the groups with more than 5 and 9 reads, respectively, would be identified.

Statistical methods

Survival analysis, including Kaplan–Meier curves plots and log-rank test *p* values, were obtained using the Python lifelines package (v0.26.4)⁹⁵. *P* values comparing TMB and aneuploidy between two groups were computed with two-sided Wilcoxon rank-sum tests. Heatmap clusters were generated through `seaborn` `clustermap`.

Viruses with significant log-rank *p* values are reported as significantly associated with survival.

None of the reference viruses were significantly associated with survival after FDR correction (Supplementary Table 1), however, we report in Fig. 2 the association between HR-HPV with unadjusted *p* value because it is confirmatory of a known association between HR-HPV and HNSC survival^{33,34}.

For HERV, our exploratory data analysis uncovered some significant associations with complete hypothesis testing. We present in the main text selected associations with at least five cases in each group. Nevertheless, FDR correction was applied within each cancer type for all HERV associations, and we additionally applied a global FDR correction for all comparisons across cancer types, yielding some significant associations with less than five positive cases. The complete significant associations between survival and viral presence are reported in Supplementary Data 12.

Reporting summary

Further information on research design is available in the Nature Portfolio Reporting Summary linked to this article.

Data availability

The complete training and test data as well as viral databases generated in this study have been deposited in the Zenodo database under the accession code <https://doi.org/10.5281/zenodo.7548375>. The results shown here are in whole or part based upon data generated by the TCGA Research Network: <https://www.cancer.gov/tcga>. The raw FASTQ RNA sequencing data are protected and are not publicly available due to data privacy laws, but are available under restricted access as data can be unique to an individual. Access can be obtained from the Genome Data Commons (GDC) after receiving permission via dbGaP, following the steps described in: https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs000178.v11.p8. The processed data including viruses identified and respective statistics are available as supplementary Data 3. The complete data generated in this study are provided in the Supplementary Information/Source Data file. Source data are provided with this paper.

Code availability

The scripts for pre and post-processing and the viRNAtrap package are available through GitHub: <https://github.com/AuslanderLab/virnatrap> and Zenodo under accession code: <https://doi.org/10.5281/zenodo.7548375>.

References

- Morales-Sánchez, A. & Fuentes-Pananá, E. M. Human viruses and cancer. *Viruses* **6**, 4047–4079 (2014).
- Krump, N. A. & You, J. Molecular mechanisms of viral oncogenesis in humans. *Nat. Rev. Microbiol.* **16**, 684–698 (2018).
- Rous, P. A sarcoma of the fowl transmissible by an agent separable from the tumor cells. *J. Exp. Med.* **13**, 397–411 (1911).
- Moore, P. S. & Chang, Y. Why do viruses cause cancer? Highlights of the first century of human tumour virology. *Nat. Rev. Cancer* **10**, 878–889 (2010).
- Salyakina, D. & Tsinoremas, N. F. Viral expression associated with gastrointestinal adenocarcinomas in TCGA high-throughput sequencing data. *Hum. Genomics* **7**, 23 (2013).
- Parfenov, M. et al. Characterization of HPV and host genome interactions in primary head and neck cancers. *Proc. Natl Acad. Sci. USA* **111**, 15544–15549 (2014).
- Cao, S. et al. Divergent viral presentation among human tumors and adjacent normal tissues. *Sci. Rep.* **6**, 28294 (2016).
- Strong, M. J. et al. Differences in gastric carcinoma micro-environment stratify according to EBV infection intensity: implications for possible immune adjuvant therapy. *PLoS Pathog.* **9**, e1003341 (2013).
- Khoury, J. D. et al. Landscape of DNA virus associations across human malignant cancers: analysis of 3,775 cases using RNA-Seq. *J. Virol.* **87**, 8916–8926 (2013).
- Tang, K. W., Alaei-Mahabadi, B., Samuelsson, T., Lindh, M. & Larsson, E. The landscape of viral expression and host gene fusion and adaptation in human cancer. *Nat. Commun.* **4**, 2513 (2013).
- Cantalupo, P. G., Katz, J. P. & Pipas, J. M. Viral sequences in human cancer. *Virology* **513**, 208–216 (2018).
- Zapatka, M. et al. The landscape of viral associations in human cancers. *Nat. Genet.* **52**, 320–330 (2020).
- Kellam, P. Molecular identification of novel viruses. *Trends Microbiol.* **6**, 160–165 (1998).
- Ren, J., Ahlgren, N. A., Lu, Y. Y., Fuhrman, J. A. & Sun, F. VirFinder: a novel k-mer based tool for identifying viral sequences from assembled metagenomic data. *Microbiome* **5**, 69 (2017).
- Ren, J. et al. Identifying viruses from metagenomic data using deep learning. *Quant. Biol.* **8**, 64–77 (2020).
- Rajkumar, U. et al. DeepViFi: detecting oncoviral infections in cancer genomes using transformers. *BCB '22: Proceedings of the 13th ACM International Conference on Bioinformatics, Computational Biology and Health Informatics 1–8* <https://doi.org/10.1145/3535508.3545551> (Association for Computing Machinery 2022).
- Fang, Z. et al. PPR-Meta: a tool for identifying phages and plasmids from metagenomic fragments using deep learning. *Gigascience* <https://doi.org/10.1093/gigascience/giz066> (2019).
- Auslander, N., Gussow, A. B., Benler, S., Wolf, Y. I. & Koonin, E. V. Seeker: alignment-free identification of bacteriophage genomes by deep learning. *Nucleic Acids Res.* **48**, e121 (2020).
- Zhang, Z. et al. Rapid identification of human-infecting viruses. *Transbound. Emerg. Dis.* **66**, 2517–2522 (2019).
- Bartoszewicz, J. M., Seidel, A. & Renard, B. Y. Interpretable detection of novel human viruses from genome sequencing data. *NAR Genom. Bioinform.* **3**, lqab004 (2021).
- Tampuu, A., Bzhalava, Z., Dillner, J. & Vicente, R. ViraMiner: deep learning on raw DNA sequences for identifying viral genomes in human samples. *PLoS ONE* **14**, e0222271 (2019).
- Weinstein, J. N. et al. The cancer genome atlas Pan-cancer analysis project. *Nat. Genet.* **45**, 1113–1120 (2013).
- Mukherjee, S., Huntemann, M., Ivanova, N., Kyrpides, N. C. & Pati, A. Large-scale contamination of microbial isolate genomes by Illumina PhiX control. *Stand. Genom. Sci.* **10**, 18 (2015).
- Altschul, S. F. et al. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res.* **25**, 3389–3402 (1997).
- Coursey, T. L., Van Doorslaer, K. & McBride, A. A. Regulation of human papillomavirus 18 genome replication, establishment, and persistence by sequences in the viral upstream regulatory region. *J. Virol.* **95**, e0068621 (2021).
- Doorbar, J. et al. The biology and life-cycle of human papillomaviruses. *Vaccine* **30**, F55–F70 (2012).
- Schäffer, A. A. et al. Integration of adeno-associated virus (AAV) into the genomes of most Thai and Mongolian liver cancer patients does not induce oncogenesis. *BMC Genomics* **22**, 814 (2021).
- Bayard, Q. et al. Cyclin A2/E1 activation defines a hepatocellular carcinoma subclass with a rearrangement signature of replication stress. *Nat. Commun.* **9**, 5235 (2018).
- Cossart, Y. E., Field, A. M., Cant, B. & Widdows, D. Parvovirus-like particles in human sera. *Lancet* **1**, 72–73 (1975).
- Adamson-Small, L. A., Ignatovich, I. V., Laemmerhirt, M. G. & Hobbs, J. A. Persistent parvovirus B19 infection in non-erythroid tissues: possible role in the inflammatory and disease process. *Virus Res.* **190**, 8–16 (2014).
- Dickinson, A. et al. Newly detected DNA viruses in juvenile nasopharyngeal angiofibroma (JNA) and oral and oropharyngeal squamous cell carcinoma (OSCC/OPSCC). *Eur. Arch. Otorhinolaryngol.* **276**, 613–617 (2019).
- Li, Y. et al. Detection of parvovirus B19 nucleic acids and expression of viral VP1/VP2 antigen in human colon carcinoma. *Am. J. Gastroenterol.* **102**, 1489–1498 (2007).
- Sethi, S. et al. Characteristics and survival of head and neck cancer by HPV status: a cancer registry-based study. *Int. J. Cancer* **131**, 1179–1186 (2012).
- Sarkar, S. et al. Human papilloma virus (HPV) infection leads to the development of head and neck lesions but offers better prognosis in malignant Indian patients. *Med. Microbiol. Immunol.* **206**, 267–276 (2017).
- Curry, G. et al. Human endogenous retrovirus K in cancer: a potential biomarker and immunotherapeutic target. *Viruses* <https://doi.org/10.3390/v12070726> (2020).

36. Kolbe, A. R. et al. Human endogenous retrovirus expression is associated with head and neck cancer and differential survival. *Viruses* <https://doi.org/10.3390/v12090956> (2020).
37. Kämmerer, U., Germeyer, A., Stengel, S., Kapp, M. & Denner, J. Human endogenous retrovirus K (HERV-K) is expressed in villous and extravillous cytotrophoblast cells of the human placenta. *J. Reprod. Immunol.* **91**, 1–8 (2011).
38. Armbruster, V. et al. A novel gene from the human endogenous retrovirus K expressed in transformed cells. *Clin. Cancer Res.* **8**, 1800–1807 (2002).
39. Wang-Johanning, F. et al. Human endogenous retrovirus K triggers an antigen-specific immune response in breast cancer patients. *Cancer Res.* **68**, 5869–5877 (2008).
40. Wang-Johanning, F. et al. Expression of human endogenous retrovirus k envelope transcripts in human breast cancer. *Clin. Cancer Res.* **7**, 1553–1560 (2001).
41. Kassiotis, G. Endogenous retroviruses and the development of cancer. *J. Immunol.* **192**, 1343–1349 (2014).
42. Xue, B., Sechi, L. A. & Kelvin, D. J. Human endogenous retrovirus K (HML-2) in health and disease. *Front. Microbiol.* **11**, 1690 (2020).
43. Kim, J. S., Yoon, S. J., Park, Y. J., Kim, S. Y. & Ryu, C. M. Crossing the kingdom border: human diseases caused by plant pathogens. *Environ. Microbiol.* **22**, 2485–2495 (2020).
44. Hahn, S. et al. Serological response to human endogenous retrovirus K in melanoma patients correlates with survival probability. *AIDS Res. Hum. Retroviruses* **24**, 717–723 (2008).
45. Zhao, J. et al. Expression of human endogenous retrovirus type K envelope protein is a novel candidate prognostic marker for human breast cancer. *Genes Cancer* **2**, 914–922 (2011).
46. Reis, B. S. et al. Prostate cancer progression correlates with increased humoral immune response to a human endogenous retrovirus GAG protein. *Clin. Cancer Res.* **19**, 6112–6125 (2013).
47. Fan, X., Quezada, S. A., Sepulveda, M. A., Sharma, P. & Allison, J. P. Engagement of the ICOS pathway markedly enhances efficacy of CTLA-4 blockade in cancer immunotherapy. *J. Exp. Med.* **211**, 715–725 (2014).
48. Xiao, Z., Mayer, A. T., Nobashi, T. W. & Gambhir, S. S. ICOS is an indicator of T-cell-mediated response to cancer immunotherapy. *Cancer Res.* **80**, 3023–3032 (2020).
49. Faget, J. et al. ICOS-ligand expression on plasmacytoid dendritic cells supports breast cancer progression by promoting the accumulation of immunosuppressive CD4+ T cells. *Cancer Res.* **72**, 6130–6141 (2012).
50. Conrad, C. et al. Plasmacytoid dendritic cells promote immunosuppression in ovarian cancer via ICOS costimulation of Foxp3(+) T-regulatory cells. *Cancer Res.* **72**, 5240–5249 (2012).
51. Zhang, T. et al. RNA viral community in human feces: prevalence of plant pathogenic viruses. *PLoS Biol.* **4**, e3 (2006).
52. Balique, F., Lecoq, H., Raoult, D. & Colson, P. Can plant viruses cross the kingdom border and be pathogenic to humans. *Viruses* **7**, 2074–2098 (2015).
53. Abbas, A. A. et al. Redondoviridae, a family of small, circular DNA viruses of the human oro-respiratory tract associated with periodontitis and critical illness. *Cell Host Microbe* **25**, 719–729.e714 (2019).
54. Halary, S. et al. Novel single-stranded DNA circular viruses in pericardial fluid of patient with recurrent pericarditis. *Emerg. Infect. Dis.* **22**, 1839–1841 (2016).
55. Dalldorf, G. & Sickles, G. M. An unidentified, filtrable agent isolated from the feces of children with paralysis. *Science* **108**, 61–62 (1948).
56. Federici, B. A. Isolation of an iridovirus from two terrestrial isopods, the pill bug, *Armadillidium vulgare*, and the sow bug, *Porcellio dilatatus*. *J. Invertebr. Pathol.* **36**, 373–381 (1980).
57. Williams, T. Natural invertebrate hosts of iridoviruses (Iridoviridae). *Neotrop. Entomol.* **37**, 615–632 (2008).
58. Li, L. et al. Investigation on mosquito-borne viruses at lancang river and Nu river watersheds in Southwestern China. *Vector Borne Zoonotic Dis.* **17**, 804–812 (2017).
59. Ohba, M. & Aizawa, K. Mammalian toxicity of an insect iridovirus. *Acta Virol.* **26**, 165–168 (1982).
60. Ince İ A. et al. Invertebrate iridoviruses: a glance over the last decade. *Viruses* <https://doi.org/10.3390/v10040161> (2018).
61. Newman, A. M. et al. Robust enumeration of cell subsets from tissue expression profiles. *Nat. Methods* **12**, 453–457 (2015).
62. Carlton, J. M. et al. Draft genome sequence of the sexually transmitted pathogen *Trichomonas vaginalis*. *Science* **315**, 207–212 (2007).
63. Kissinger, P. *Trichomonas vaginalis*: a review of epidemiologic, clinical and treatment issues. *BMC Infect. Dis.* **15**, 307 (2015).
64. Yang, S. et al. *Trichomonas vaginalis* infection-associated risk of cervical cancer: a meta-analysis. *Eur. J. Obstet. Gynecol. Reprod. Biol.* **228**, 166–173 (2018).
65. Risinger, J. I. et al. PTEN mutation in endometrial cancers is associated with favorable clinical and pathologic characteristics. *Clin. Cancer Res.* **4**, 3005–3010 (1998).
66. Barretina, J. et al. The cancer cell line encyclopedia enables predictive modelling of anticancer drug sensitivity. *Nature* **483**, 603–607 (2012).
67. Banerjee, S. et al. The ovarian cancer oncobiome. *Oncotarget* **8**, 36225–36245 (2017).
68. Nejman, D. et al. The human tumor microbiome is composed of tumor type-specific intracellular bacteria. *Science* **368**, 973–980 (2020).
69. Robinson, H. L. Retroviruses and cancer. *Rev. Infect. Dis.* **4**, 1015–1025 (1982).
70. Uphoff, C. C., Lange, S., Denkmann, S. A., Garritsen, H. S. & Drexler, H. G. Prevalence and characterization of murine leukemia virus contamination in human cell lines. *PLoS ONE* **10**, e0125622 (2015).
71. Kotic, A. D. et al. PathSeq: software to identify or discover microbes by deep sequencing of human tissue. *Nat. Biotechnol.* **29**, 393–396 (2011).
72. Ahlers, L. R., Bastos, R. G., Hiroyasu, A. & Goodman, A. G. Invertebrate iridescent virus 6, a DNA virus, stimulates a mammalian innate immune response through RIG-I-Like receptors. *PLoS ONE* **11**, e0166088 (2016).
73. Twu, O. et al. *Trichomonas vaginalis* exosomes deliver cargo to host cells and mediate host: parasite interactions. *PLoS Pathog.* **9**, e1003482 (2013).
74. Wu, X. et al. Identification of key genes and pathways in cervical cancer by bioinformatics analysis. *Int. J. Med. Sci.* **16**, 800–812 (2019).
75. Taylor, L. J. et al. Redondovirus diversity and evolution on global, individual, and molecular scales. *J. Virol.* **95**, e0081721 (2021).
76. Hatcher, E. L. et al. Virus variation resource - improved response to emergent viral outbreaks. *Nucleic Acids Res.* **45**, D482–D490 (2017).
77. Sayers, E. W. et al. Database resources of the National Center for Biotechnology Information. *Nucleic Acids Res.* **49**, D10–D17 (2021).
78. Chollet, F. et al. Keras. <https://github.com/fchollet/keras> (2015).
79. Grossman, R. L. et al. Toward a shared vision for cancer genomic data. *N. Engl. J. Med.* **375**, 1109–1112 (2016).
80. Van Doorslaer, K. et al. The papillomavirus episteme: a major update to the papillomavirus sequence database. *Nucleic Acids Res.* **45**, D499–D506 (2017).
81. Goodacre, N., Aljanahi, A., Nandakumar, S., Mikailov, M. & Khan, A. S. A reference viral database (RVDB) to enhance bioinformatics analysis of high-throughput sequencing for novel virus detection. *mSphere* <https://doi.org/10.1128/mSphereDirect.00069-18> (2018).

82. Tokuyama, M. et al. ERVmap analysis reveals genome-wide transcription of human endogenous retroviruses. *Proc. Natl Acad. Sci. USA* **115**, 12565–12572 (2018).
83. Paces, J. et al. HERVd: the human endogenous retroviruses database: update. *Nucleic Acids Res.* **32**, D50 (2004).
84. Karolchik, D. et al. The UCSC Table Browser data retrieval tool. *Nucleic Acids Res.* **32**, D493–D496 (2004).
85. Yutin, N., Puigbò, P., Koonin, E. V. & Wolf, Y. I. Phylogenomics of prokaryotic ribosomal proteins. *PLoS ONE* **7**, e36972 (2012).
86. Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Mol. Biol. Evol.* **30**, 772–780 (2013).
87. Bannert, N. & Kurth, R. Retroelements and the human genome: new perspectives on an old relation. *Proc. Natl Acad. Sci. USA* **101**, 14572–14579 (2004).
88. Smith, C. C. et al. Endogenous retroviral signatures predict immunotherapy response in clear cell renal cell carcinoma. *J. Clin. Invest.* **128**, 4804–4820 (2018).
89. Dobin, A. et al. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15–21 (2013).
90. Schäffer, A. A. et al. VecScreen_plus_taxonomy: imposing a tax(onomy) increase on vector contamination screening. *Bioinformatics* **34**, 755–759 (2018).
91. Wood, D. E., Lu, J. & Langmead, B. Improved metagenomic analysis with Kraken 2. *Genome Biol.* **20**, 257 (2019).
92. Celaj, A., Markle, J., Danska, J. & Parkinson, J. Comparison of assembly algorithms for improving rate of metatranscriptomic functional annotation. *Microbiome* **2**, 39 (2014).
93. Taylor, A. M. et al. Genomic and functional approaches to understanding cancer aneuploidy. *Cancer Cell* **33**, 676–689.e673 (2018).
94. Goldman, M. J. et al. Visualizing and interpreting cancer genomics data via the Xena platform. *Nat. Biotechnol.* **38**, 675–678 (2020).
95. Davidson-Pilon, C. Lifelines: survival analysis in Python. *J. Open Source Softw.* **4**, 1317 (2019).

Acknowledgements

The research reported in this publication was supported in part by the National Cancer Institute of the National Institutes of Health under Award Number R00CA252025 (N.A.), RF1-AG063481, P30-CA016520 (Q.L.), and NIH RO1 AI153508, Commonwealth of Pennsylvania SAP# 4100089371, and P30 CA010815 (P.M.L.), and by the Intramural Research Program of the National Institutes of Health, National Cancer Institute (A.A.S.).

Author contributions

N.A. initiated the project. Q.L., R.Z., A.A.S., and N.A. supervised work. A.E., X.H., R.Z., A.A.S., and N.A. designed and performed experiments and analyses. A.E., Y.Y., D.E.S., J.W., A.A.S., and N.A. wrote and tested software. P.M.L. and Q.M. contributed to data interpretation and exploratory analyses. K.T. and Q.L. revised the survival analysis.

Competing interests

P.M.L. is a founder of and advisor to Vironika, LLC. All other authors declare no competing interests.

Additional information

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s41467-023-36336-z>.

Correspondence and requests for materials should be addressed to Noam Auslander.

Peer review information *Nature Communications* thanks the anonymous reviewers for their contribution to the peer review of this work. Peer reviewer reports are available.

Reprints and permissions information is available at <http://www.nature.com/reprints>

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2023