# The landscape of viral associations in human cancers

Marc Zapatka[1]*, Ivan Borozan[2]*, Daniel S. Brewer[3,4]*, Murat Iskar[1]*, Adam Grundhoff[5], Malik Alawi[5,6], Nikita Desai[7,8], Holger Sültmann[9,10], Holger Moch[11], PCAWG-Pathogens[12], Colin S. Cooper[4,13], Roland Eils[14,15,16], Vincent Ferretti[17,18], Peter Lichter[1,10], PCAWG Consortium[19]

1 Division of Molecular Genetics, German Cancer Research Center (DKFZ), Heidelberg, Germany.

2 Informatics and Bio-computing Program, Ontario Institute for Cancer Research, Toronto, Ontario, Canada

3 Norwich Medical School, University of East Anglia, Norwich, UK

4 Earlham Institute, Norwich, UK

5 Heinrich-Pette-Institute, Leibniz Institute for Experimental Virology, Hamburg, Germany

6 Bioinformatics Core, University Medical Center Hamburg-Eppendorf, Hamburg, Germany

7 Bioinformatics Group, Department of Computer Science, University College London, London, UK

8 Biomedical Data Science Laboratory, Francis Crick Institute, London UK

9 Division of Cancer Genome Research, German Cancer Research Center (DKFZ) and National Center for Tumor Diseases (NCT), Heidelberg, Germany

10 German Cancer Consortium (DKTK), Heidelberg, Germany

11 Department of Pathology and Molecular Pathology, University and University Hospital Zürich, Zürich, Switzerland

12 A full list of authors can be found at the end of the article

13 The Institute of Cancer Research, London, UK

14 Division of Theoretical Bioinformatics, German Cancer Research Center (DKFZ), Heidelberg, Germany

15 Department of Bioinformatics and Functional Genomics, Institute of Pharmacy and Molecular Biotechnology, Heidelberg University and BioQuant Center, Heidelberg, Germany

16 Center for Digital Health, Berlin Institute of Health and Charité Universitätsmedizin Berlin, Berlin, Germany

17 Ontario Institute for Cancer Research, MaRS Centre, Toronto, Canada

18 Department of Biochemistry and Molecular Medicine, University of Montreal, Montreal, Canada.

19 A full list of authors can be found in the Supplementary Note

## Equal contributions statement

MZ, IB, DSB and MI contributed equally.

## Corresponding author statement

Peter Lichter peter.lichter@dkfz-heidelberg.de

# 36Abstract

37Here, as part of the Pan-Cancer-Analysis-of-Whole-Genomes (PCAWG), which aggregated 38whole genome and, for a subset, transcriptome sequencing data from 2,658 cancers across 38 39tumor types, we systematically investigated potential viral pathogens using a consensus 40approach integrating three independent pipelines. Viruses were detected in 382 genome and 4168 transcriptome datasets. We showed the high prevalence of known tumor-associated-42viruses such as EBV, HBV and HPV16/18. The study revealed significant exclusivity of HPV 43with driver mutations in head-and-neck cancer and associated HPV with APOBEC 44mutational signatures, suggesting a role of impaired antiviral defense as driving force in 45cervical, bladder and head-and-neck carcinoma. For HBV, HPV16/18 and AAV2 viral 46integration was associated with local variations in genomic copy number. Integrations at 47the *TERT* promoter were coupled to high telomerase expression evidently activating this 48tumor driving process. High levels of endogenous retrovirus ERV1 expression were linked to 49worse survival outcome in kidney cancer.

# 50Introduction

51The World Health Organization estimates that 15.4% of all cancers are attributable to 52infections and 9.9% are linked to viruses[1,2]. Cancers attributable to infections have a greater 53incidence than any individual type of cancer worldwide. Eleven pathogens have been 54classified as carcinogenic agents in humans by the International Agency for Research on 55Cancer (IARC)[3]. After *Helicobacter pylori* (associated with 770,000 cases), the four most 56prominent infection related causes of cancer are estimated to be viral[2]: human papilloma virus 57(HPV)[4,5] (associated with 640,000 cancers), hepatitis B virus (HBV)[5] (420,000), hepatitis C 58virus (HCV)[6] (170,000) and Epstein-Barr Virus (EBV)[7] (120,000). It has been shown that 59viruses can contribute to the biology of multistep oncogenesis and are implicated in many of 60the hallmarks of cancer[8]. Most importantly, the discovery of links between infection and 61cancer types has provided actionable opportunities, such as HPV vaccines as preventive 62measure, to reduce the global impact of cancer. The following characteristics were proposed 63to define human viruses causing cancer through direct or indirect carcinogenesis[9]: i) Presence 64and persistence of viral DNA in tumor biopsies; ii) Growth promoting activity of viral genes 65in model systems; iii) Dependence of malignant phenotype on continuous viral oncogene 66expression or modification of host genes; iv) Epidemiological evidence that a virus infection 67represents a major risk for development of cancer.
68
69The worldwide efforts of comprehensive genome and transcriptome analyses of tissue 70samples from cancer patients generate appropriate facilities for capturing information not 71only from human cells, but also from other - potentially pathogenic - organisms or viruses 72present in the tissue. A comprehensive collection of whole genome and transcriptome data 73from cancer tissues has been generated within the ICGC (International Cancer Genome 74Consortium) project PCAWG (Pan-Cancer Analysis of Whole Genomes)[10], providing a 75unique opportunity for a systematic search for tumor-associated viruses.
76
77The PCAWG Consortium aggregated whole genome sequencing data from 2,658 cancers 78across 38 tumor types generated by the ICGC and TCGA projects. These sequencing data 79were re-analyzed with standardized, high-accuracy pipelines to align to the human genome 80(build hs37d5) and identify germline variants and somatically acquired mutations[10]. The 81PCAWG working group "Exploratory Pathogens" analyzed the whole genome sequencing 82(WGS) and whole transcriptome sequencing (RNA-seq) data of the PCAWG consensus

83cohort (2,656 donors). Focusing on viral pathogens, we applied three independently 84developed pathogen detection pipelines 'Computational Pathogen Sequence Identification' 85(CaPSID)[11], 'Pathogen Discovery Pipeline' (P-DiP), and 'SEarching for PATHogens' 86(SEPATH) to generate a large compendium of viral associations across 38 cancer types. We 87extensively characterized the known and novel viral associations by integrating driver 88mutations, mutational signatures, gene expression profiles and patient survival data of the 89same set of tumors analyzed in PCAWG.

# 90Results

## 91Identification of tumor-associated viruses

92To identify the presence of viral sequences, we explored the WGS data of 5,354 93tumor/normal samples across 38 cancer types, and 1,057 tumor RNA-seq data across 25 94cancer types (Supplementary Tables 1,2,20). 195.8 billion reads were considered for analysis 95as they were not sufficiently aligned to the human reference genome in the PCAWG-96generated alignment. Remaining reads ranged from 28,036 to 800 million per WGS and up to 97120 million per RNA-seq tumor sample (Fig. 1a, Extended Data Figure 1a-c). Viral 98sequences were detected and quantified independently by three recently developed pathogen 99discovery pipelines CaPSID, P-DiP and SEPATH. The estimated relative abundance of a 100virus was calculated as viral reads per million extracted reads (PMER) at the genus level to 101improve consistency between pipelines. To minimize the rate of false positives in virus 102detection, we applied a strict threshold of PMER>1 supported by at least three viral reads as 103similarly suggested by previous studies[11,12]. Virus detection in a sample by at least two 104pipelines was considered as a consensus hit. In total, 532 genera were considered for the 105extensive virus search in at least two of the pipelines (Extended Data Figure 1d, 106Supplementary Table 18). Filtering of suspected viral laboratory contaminants was achieved 107through P-DiP, by examining each assembled contig of viral sequence segments for artificial, 108non-viral vector sequences and inspecting virus genome coverage across all positive samples 109(Extended Data Figure 2a). The most frequent hits prone to suspected contamination were 110lambdavirus, alphabaculovirus, microvirus, simplexvirus, hepacivirus, cytomegalovirus, 111orthopoxvirus and punalikevirus; these were observed across many tumor types (Fig. 1b). For 112example, mastadenovirus showed an uneven genome coverage which could result from 113contaminating vector sequences. Therefore, we analyzed the virus detections across 114sequencing dates (Extended Data Figure 2b) to assess any batch effect indicative of a 115contaminant; in mastadenovirus, we identified an association with sequencing date in early-116onset prostate cancer regardless of tumor/normal state. We conclude that our mastadenovirus 117detections are due to a contamination occurring across projects worldwide where similar 118patterns could be identified.
119
120We generally observed a strong overlap of the genera identified across pipelines (Extended 121Data Figure 1e, Supplementary Tables 6,7,11). From the whole genome dataset, we identified 122321, 598 and 206 virus-tumor pairs for P-DiP, CaPSID and SEPATH, respectively (Fig. 2a, 123overlap after random permutation of detections, Extended Data Figure 3a, Supplementary 124Tables 3-5). The number of hits derived from the RNA-seq dataset differed between the 125pipelines (virus-tumor pairs: 101 for P-DiP, 83 for CaPSID, 41 for SEPATH; Fig. 2b, 126Supplementary Tables 8-10). SEPATH, using a k-mer approach, detected the lowest number 127of virus hits and was the least sensitive. Despite this, the identified viruses matched well with 128the consensus (DNA 90%, RNA 95%). P-DiP, based on an assembly and BLAST approach, 129detected more hits with 59% of the DNA and 54% of the RNA hits in the consensus set, 130while CaPSID, being most sensitive, implementing a two-step alignment process 131complemented by an assembly step, identified 60% (DNA) and 80% (RNA) hits within the 132consensus set. While the majority of the virus hits from RNA-seq (n=61/68) were 133overlapping with the WGS data, a lower fraction of detections in WGS data were present in 134the RNA-seq data (n=61/168 of 382 virus detections with RNA-seq data), emphasizing the 135importance of DNA sequencing for generating an unbiased catalogue of tumor-associated

11
12

136viruses. This difference can also be attributed to the viral life cycle as during incubation or
137latent phases, viral gene expression can be minimal[13]. Contrasting virus positive and negative
138samples within each organ type shows that the organ system, as expected, has a significant
139influence, but not virus positivity ($P < 2 \times 10^{-16}$, ANOVA modeling candidate reads
140dependent on organ system and virus positivity, Extended Data Figure 1c). This indicates that
141virus-positive tumors were not detected due to a higher number of candidate reads and is in
142line with the fact that the viral reads in most cases do not substantially contribute to the reads
143analyzed. 86% of the sequence hits detected from WGS and RNA-seq data were found to be
144from double-stranded DNA viruses (dsDNA) and dsDNA viruses with reverse transcriptase
145(Fig. 1c, Supplementary Table 19). This could be attributed to i) a higher frequency of tumor-
146associated viruses from these genome types[15], ii) a larger sequence dataset for WGS in
147comparison to RNA-seq, iii) a potential limitation of our analysis due to DNA and RNA
148extraction protocols that are less likely to include single-stranded (ss)DNA or RNA viruses or
149iv) the selection bias of tumor entities included in the PCAWG study (Fig. 1c).

## 150The virome landscape across 38 distinct tumor types

151We employed a consensus approach that resulted in a reliable set of 389 distinct virus-tumor
152pairs from WGS and RNA-seq data (Fig. 2a-d). Overall, 23 virus genera were detected across
153356 tumor patients (13%). The top five most prevalent viruses (lymphocryptovirus,
154orthohepadnavirus, roseolovirus, alphapapillomavirus, cytomegalovirus) account for 85% of
155the consensus virus hits in tumors (n=329 out of 389). Among these five prevalent virus
156genera, three have been well described in the literature as drivers of tumor initiation and
157progression[9]: i) lymphocryptovirus (n=145 samples, 5.5%, e.g. Epstein-Barr Virus, EBV) is
158the most common viral infection across a variety of tumor entities mainly from
159gastrointestinal tract, and showed a much lower prevalence in the matched non-malignant
160control samples (n=82, 3%) (Fig. 2c); ii) orthohepadnavirus (n=67, 2.5%, e.g. hepatitis B,
161HBV) are as expected the most frequent among liver cancer with HBV present in 62 of 330
162donors (18.9%); and iii) alphapapillomavirus (discussed below). Lymphocryptovirus (n=11),
163orthohepadnavirus (n=18) and alphapapillomavirus (n=32) were detected both in RNA and
164DNA sequencing data (Fig. 2c, left panel), with alphapapillomavirus being the most frequent
165(32 out of 39 consensus hits). This is in line with the constitutive expression of viral
166oncogenes in cancers associated with these viruses, a parameter supporting a direct role in
167carcinogenesis[9]. An in-depth analysis of the virus genome equivalents per human tumor
168genome equivalent considering genome sizes, coverage and tumor purity showed overall low
169viral genome equivalents even for established tumor viruses (Extended Data Figure 3c,
170Supplementary Table 12). Evidence for mouse mammary tumor virus (MMTV, PMER = 3.4)
171was detected in one renal carcinoma sample and in none of the 214 analyzed breast cancer
172samples. Previous work has suggested that MMTV may play a role in breast cancer but our
173comprehensive search of viral sequences could not identify any MMTV-positive case in
174breast cancer that would support this claim.
175
176Roseolovirus and Alphatorquevirus show a higher number of hits in the non-malignant
177control samples, which were mainly derived from blood cells (Fig. 2c). For example, we
178identified 59 patients as Roseolovirus-positive (HHV-6A, HHV-6B, HHV-7) in their tumor
179(pancreas 6%, stomach 8%, colon/rectum 8.3%) and 90 patients positive in the non-malignant
180control samples. Considering the known cell tropism of roseolovirus for B- and T-cells[14], we
181asked whether immune infiltration would be higher in roseolovirus-positive tumors.
182However, we could not identify a stronger contribution of immune cells in virus positive
183tumor cases as estimated using CIBERSORT[15] (false discovery rate (FDR) corrected $P > 0.05$
184for pancreas; Extended Data Figure 4a). Therefore, in line with current knowledge (reviewed
185in[16]), we cannot confirm a link between roseolovirus and immune cell content or tumor
186development. Furthermore, we could not identify actively transcribed viral genes for
187Roseolovirus and Alphatorquevirus at the transcriptome level. This is in agreement with the
188latent state of these viruses reported for blood mononuclear cells[14], and their transmission
189through blood transfusions[17]. Cytomegalovirus (CMV) was found, as expected[18], after
190identifying and removing contaminations both in stomach tumors (n=13) and the adjacent

191non-malignant tissue (n=11). In line with a recent publication[19], we could not detect CMV in 192the analyzed 294 CNS-tumors (146 medulloblastomas, 89 pilocytic astrocytoma, 41 193glioblastomas, 18 oligodendrogliomas). Therefore, a previously debated role of this virus is 194not supported. Notably, we did not identify a significant enrichment of co-infection of 195multiple viruses in any tumor type (Extended Data Figure 3d).

## 196Hepatitis B virus

197Hepatitis B virus was most frequently detected among liver cancers (n=62). Comparing to the 198histopathological gold standard HBV PCR test[20,21] (n=228), we found the WGS-based 199consensus detections had the same high specificity (96.1%) and a high sensitivity (84.0%), 200indicating that the HBV detections by WGS are reliable (Fig. 3a, Extended Data Figure 4b, 201Supplementary Table 13). Furthermore, five out of seven cases positive in WGS and negative 202for HBV PCR showed positivity for HBAg indicating a high sensitivity of the WGS analysis. 203In summary, the precision (85.7%) and recall (84%) for the detection of HBV based on ~30x 204WGS is comparable to targeted PCR. We confirmed a significant exclusivity between HBV 205infection and *CTNNB1*, *TP53* and *ARID1A* mutations that was found in a larger liver cancer 206cohort analyzed by high throughput sequencing (FDR corrected $P = 5.35 \times 10^{-6}$, 0.0023 and 2070.0023, DISCOVER[22])[23].

## 208Epstein-Barr virus

209Epstein-Barr virus was detected in many different tumor entities and normal samples (Fig. 2102c). Comparing EBV PMERs in tumors and matched normals we see a stronger contribution 211in matched normals from matched solid tissue or tissue adjacent to the tumor (Extended Data 212Figure 4c). For samples showing reads for EBV in WGS and with available RNA sequencing 213data, the absolute score for immune cells based on CIBERSORT[15] was not significantly 214different between virus positive and negative samples (FDR corrected $P > 0.05$ for 215colon/rectum, head/neck, lymphoid, stomach; Extended Data Figure 4a). In summary, there is 216no evidence for a detection of EBV due to infiltrating immune cells. This indicates EBV 217presence in the respective organs. Based on the expression data available for the tumor 218samples we identified viral transcripts of the latent as well as lytic phase of the viral lifecycle 219(Fig. 3b, Extended Data Figure 4d, Supplementary Table 13). Eight of the nine tumors 220expressing lytic EBV transcripts are from stomach, confirming the active contribution of 221EBV to gastric cancer[24].
222

## 223Alphapapillomavirus

224Alphapapillomaviruses were mainly detected in head-and-neck cancer (n=18 of 57), cervical 225cancer (n=19 of 20) and in two bladder cancer cases out of 23, in agreement with previous 226studies[4,25,26]. There is also supporting evidence for 32 out of 39 alphapapillomavirus hits in the 227transcriptome data (Fig. 2c). We observed only one HPV subtype per tumor according to the 228P-DiP results with HPV16 being the dominant type in cervix (n=11) and head-and-neck 229(n=15) tumors, followed by HPV18 only present in cervical cancer (n=6). As reported 230previously[27], HPV33 was identified in head-and-neck (n=3) and cervix (n=1) tumors. 231Different HPV variants, type 6 and 45, were detected in bladder cancer.
232
233In head-and-neck cancer, HPV-positive tumors exhibit an almost complete mutual exclusivity 234with mutations in known drivers like *TP53*, *CDKN2A* and *TERT* (FDR corrected $P = 1.73 \times$ 235$10^{-5}$, $1.73 \times 10^{-5}$, 0.012; multiple testing corrected for presented mutations in EBV and HPV, 236DISCOVER[22]) (Fig. 3c, Supplementary Table 13), as reported previously[25], which could be 237explained by a mutation independent inactivation of TP53 through the human 238papillomaviruses[28–30]. Furthermore, we identified mutational signature 2 as enriched for 239alphapapillomavirus positive cases in head-and-neck cancers (FDR corrected $P=0.02$; Fig. 2403d, Supplementary Table 12,22)[31]. In addition, the expression of APOBEC3B is significantly

18
19
20

241higher in the virus positive head-and-neck cancers compared to their negative counterparts
242($P$=1.6 × 10$^{-4}$, Fig. 3f)[32]. However, we did not observe enrichment of APOBEC signatures
243and expression changes for EBV-positive samples either in cervix or in other tissues.
244
245Distinct expression profiles between virus positive and negative tumors in head-and-neck
246cancer are observed (Fig. 3e, Supplementary Table 23)[33]. Analyzing the immune cells
247estimated by CIBERSORT, we identified a significant increase in macrophages and T-cell
248signals in alphapapillomavirus positive head-and-neck cancers ($P$=0.004, 0.012 and 0.012 for
249follicular helper, CD8 and regulatory T-cells and $P$=0.018 for M1-Macrophages; FDR
250corrected for all viruses and cell types tested; Fig. 3g, Supplementary Table 24). Our
251integrative analysis on HPV reconfirms many of the findings related to HPV infection,
252illustrating the potential of our systematic approach in identifying and characterizing tumor-
253associated viruses.

## 254Activation of endogenous retroviruses linked to outcome

255Human endogenous retroviruses (HERV) are integrated in the human DNA originating from
256infection of germline cells by retroviruses over millions of years[34] and contribute over
257500,000 individual sites, or 2.7% of the overall sequence the human genome[35,36]. The
258endogenous retroviruses were identified by the three pathogen detection pipelines but filtered
259by CaPSID and SEPATH. In addition, an alignment-based approach was used to detect
260HERV sequences embedded in the human reference genome that could be missed by the
261pipelines focusing only on non-human reads. In this study, we quantified the expression of
262HERV-like LTR (long terminal repeat) retrotransposons categorized into several clades by
263Repbase[37] as ERVL, ERVL-MaLR, ERV1, ERVK and ERV (Supplementary Table 14). In
264comparison to the other HERV families, ERV1 shows the strongest expression on average
265(Fig. 4a) and ERVK the highest fraction of active loci (Fig. 4b). Analyzing the expression of
266HERVs we could identify a strong expression for ERV1 in chronic lymphocytic leukemia
267compared to all other tumor tissues and adjacent normal tissues (Fig. 4c). However, we could
268not identify a link between transcriptionally active stemness markers (OCT3/4, SOX2, KLF4)
269and increased HERV expression, in contrast to what was reported in Ohnuki et al.[38]
270(Spearman Rank correlation < 0.35, Extended Data Figure 5). New data suggest that
271expression of HERVs is associated with prognosis in clear cell renal cell carcinoma
272(ccRCC)[39]. Analyzing the HERV expression in relation to patient survival, we identified a
273high ERV1 expression in kidney cancer linked to worse survival outcome ($P$=0.0081; Log-
274rank test; Fig. 4d, Extended Data Figure 6, Supplementary Table 15).

## 275Genomic integration of viral sequences

276Viral integration into the host genome has been shown to be a causal mechanism that can lead
277to cancer development[40]. This process is well-established for human papilloma viruses
278(HPVs) in cervical, head-and-neck and several other carcinomas, and for hepatitis B virus
279(HBV) in liver cancer[41,42].
280
281Low confidence integration events were detected for the HHV4 (gastric cancer and malignant
282lymphoma) and HPV6b (head-and-neck and bladder carcinoma), while integration events
283with high confidence were demonstrated for HBV (liver cancer), Adeno-associated virus-2
284(AAV2) (liver), HPV16/18 (both in cervical and head-and-neck carcinoma). Most of these
285integration events were found to be distributed across chromosomes and a significant number
286of viral integrations occur in the intronic (40%) regions while only 3.4% were detected in
287gene coding regions (Extended Data Figure 7a-d).
288
289HBV was found to be integrated in 36 liver cancer specimens out of 61 patients identified as
290HBV-positive. Notably, genomic clusters of viral integrations were identified in *TERT*
291(ngc=6, number of integration sites within a genomic cluster), *KMT2B* (ngc=4), recently
292identified to be a likely cancer driver gene[43,44] and *RGS12* (ngc=3)(Extended Data Figure 7e).

293Furthermore, two or more integration events in individual samples were observed in the gene
294(or gene promoter) regions of *CCCNE1*, *CDK15*, *FSIP2*, *HEATR6*, *LINC01158*, *MARS2* and
295*SLC1A7* (Fig. 5a). Additional events with two integration sites were also detected within a 50
296kb distance away from *CLMP*, *CNTNAP2* and *LINC00359* genes. Integration events at *TERT*
297were found to recur in five different liver cancer samples. One sample had a genomic cluster
298of three viral integration events within *TERT* and four samples contained a single integration
299event in the *TERT* promoter, (3') or 5' UTR regions (Supplementary Table 17). When
300comparing gene expression in samples with virus integration to those without, only TERT
301was over-expressed (fold change ≥ 2.0) in two liver cancer samples (Fig. 5e). Additional
302genes with increased expression impacted by integration events include *TEKT3*, *CCNA2*,
303*CDK15* and *THRB* (Fig. 5a).
304

305There was a significant association between HBV viral integrations and somatic copy number
306alterations (SCNAs, Fig. 5c). For samples with HBV integration events, the number of
307SCNAs was higher on average in the vicinity of viral integration sites (within 1 Mb) when
308compared to samples without HBV integration (mean: 4.2 vs 2.3, $P=7.4 \times 10^{-3}$; two-sided
309paired *t*-test). No evidence for an SCNA association was seen for other integrated viruses like
310HPV16/18 (Extended Data Figure 8a-b).
311

312HPV18 integration events were detected in seven tumors in total (Fig. 5b), with the most
313notable clusters of integration events in cervical cancer samples affecting *TALDO1* (ngc = 4)
314(Extended Data Figure 7g).
315

316In 20 samples, HPV16 integration events were detected. Genomic clusters of viral integration
317sites were identified in cervical and head-and-neck cancer samples (Extended Data Figure
3187f). None of these multiple integration events were observed to recur across patients (Fig.
3195b). Integration events were also observed in two different lncRNAs, *LINC00111* and the
320plasmacytoma variant translocation 1 gene (*PVT1*), an oncogenic lncRNA[45,46]. Expression of
321both genes is strongly increased in the cases with HPV16 integration (Extended Data Figure
3228f, Supplementary Table 17).
323

324Using the PCAWG single nucleotide variant (SNV) calls[10] we have found a significant
325increase in the number of mutations occurring within +/- 10,000 bp of high-confidence viral
326integration sites (average number of mutations per sample = 0.41 (HPV16+) vs 0.14
327(HPV16-), $P=0.02$; paired *t*-test one-sided, alternative greater, Extended Data Figure 8cd).
328Interestingly the integration sites are, compared to a random genome background, enriched in
329close proximity (<1000 bp) to common fragile sites ($P=0.0018$, Kolmogorov–Smirnov test).
330These results suggest that HPV16 integration reflects either characteristics of chromatin
331features that favor viral integration, such as fragile sites or regions with limited access to
332DNA repair complexes, or the influence of integrated HPV16 on the host genome. Such a
333correlation was not seen for the integration sites of other viruses (Extended Data Figure 8e).
334Finally, a single AAV2 integration event located in the intronic region of the cancer driver
335gene *KMT2B*[47] was detected in one liver cancer sample.

## 336Identification of novel viral species or strains

337De novo analysis using the CaPSID-pipeline has generated 56 different contigs that have
338been classified into taxonomic groups at the genus level by CSSSCL[48]. After filtering de novo
339contigs for their homology to known reference sequences, we have identified 29 contigs in 28
340different tumor samples showing low sequence similarity (in average 63%) to any nucleotide
341sequence contained in the BLAST database. In this respect, our analysis has shown that WGS
342and RNA-seq can be used to identify isolates from potentially new viral species. However,
343the total number of novel isolates were quite low in comparison to viral hits to well-defined
344genera (Fig. 2c). These *de novo* contigs were not enriched for a specific tumor entity but
345rather distributed across cancer types including bladder, head/neck and cervical cancers
346(Extended Data Figure 9).

# 347Discussion

348Searching large pan-cancer genome and transcriptome data sets allowed the identification of 349an unexpectedly high percentage of virus associated cases (16%). In particular, analysis of 350tumor genomes, which were sequenced on average to a depth of at least 30-fold coverage, 351identified considerably more virus positive cases than investigations of transcriptome data 352alone, which is the search space analyzed in most previous virome studies. This is probably 353mainly due to viruses with no or only weak transcriptional activity in the given tumor tissue. 354Co-infections, generally believed to indicate a weak immune system, were very rare 355(Extended Data Figure 3d). This could, however, also be the result of selection processes 356during tumorigenesis.
357

358While universal criteria for a causality of viral pathogens are prone to errors, it is worthwhile 359to look at individual features that might support a potentially pathomechanistic contribution 360of a given pathogen. These include aspects that affect the expression of host factors, e.g. upon 361viral integration, or the mutual exclusivity of the presence of viral genomes and other host 362factors, which are already known to play a role in the etiology of a given tumor type. Such 363aspects need to be carefully considered when discussing of what strengthens a potentially 364pathogenic role of virus.
365

366Not surprisingly, known tumor associated viruses, such as EBV, HBV and HPV16/18, were 367among the most frequently detected targets. Interestingly, viral detection based on whole 368genome sequencing showed similar performance with respect to precision and recall as a 369targeted PCR for HBV indicating the sensitivity of this approach to detect viruses. This is 370particularly true for the common integration verified for HBV and HPV 16/18 in our study. In 371addition, the common theme of potential pathomechanistic effects by the genomic integration 372of viruses, also supported by the observations of multiple nearby integration sites in a given 373tumor genome that we also report in the present study, has gained further momentum. 374Analyzing the effect of viral integrations on gene expression, we identified several links to 375genes nearby the integration site. In this regard, the frequently observed integration of HBV 376at the *TERT* promoter accompanied with the transcriptional upregulation of *TERT*, constitutes 377an intriguing mechanistic example, since an increased activity of TERT is a well-understood 378driver of carcinogenesis[49]. Furthermore, we also linked viral integrations to increased 379mutations (SNVs and SCNAs) nearby the integration site.
380

381The known causal role of HPV16/18 in several tumor entities, that triggered one of the largest 382measures in cancer prevention, has been the motivation for extensive elucidation of the 383pathogenetic processes involved. Nevertheless, comprehensive analyses of WGS and RNA-384seq data sets revealed additional novel findings. While we confirmed the exclusivity of HPV 385infection and *TP53*, *CDKN2A* and *TERT* mutations in head-and-neck tumors, we could also 386link virus presence to an increase in mutations attributed to the mutational signature 2[50]. 387These are explained by the activity of APOBEC, which – among other effects – changes viral 388genome sequences as a mechanism of cellular defense against viruses[51,52]. This activation 389could play an important role in introducing further host genome alterations and, thus, 390constitute an important mechanism driving tumorigenesis[32,52]. In liver cancer mutations in 391*CTNNB1*, *TP53* and *ARID1A*, major primary oncogenes in this cancer type and HBV 392infections were confirmed to occur significantly exclusive[23]. Furthermore, the virus positive 393head-and-neck cancer samples had a significantly higher abundance of T-cell and M1 394macrophage expression signals, which matches with the recently described subtypes of 395HNSCC that differ – among others – in virus infection and inflammation features.
396

397

## 398Acknowledgements

417

## 418Author Contributions

419MI, DB, IB, MZ contributed equally, MZ, PL jointly supervised research. VF, RE, CC, MI, 420IB, MZ, PL conceived and designed the experiments. HS performed experiments. MI, DB, 421IB, MZ performed statistical analysis. ND, MI, AG, DB, IB, MZ analysed the data. VF, RE, 422CC, HM, MA, AG, DB, IB, MZ contributed reagents/materials/analysis tools. MI, DB, IB, 423MZ, PL wrote the paper. VF, AG, CC, DB, MI, IB, MZ and PL critiqued manuscript for 424intellectual content.

## 425Competing Interests Statement

426The authors have declared that they have no competing interests.

## 427References

4281. Parkin, D. M. The global health burden of infection-associated cancers in the year 4292002. Int. J. cancer 118, 3030–44 (2006).
4302. Plummer, M. et al. Global burden of cancers attributable to infections in 2012: a 431synthetic analysis. Lancet. Glob. Heal. 4, e609-16 (2016).
4323. Bouvard, V. et al. A review of human carcinogens—Part B: biological agents. Lancet 433Oncol. (2009). doi:10.1016/S1470-2045(09)70096-8
4344. Muñoz, N., Castellsagué, X., de González, A. B. & Gissmann, L. Chapter 1: HPV in 435the etiology of human cancer. Vaccine 24 Suppl 3, S3/1-10 (2006).
4365. Bialecki, E. S. & Di Bisceglie, A. M. Clinical presentation and natural course of 437hepatocellular carcinoma. Eur. J. Gastroenterol. Hepatol. 17, 485–9 (2005).
4386. Hermine, O. et al. Regression of splenic lymphoma with villous lymphocytes after 439treatment of hepatitis C virus infection. N. Engl. J. Med. 347, 89–94 (2002).
4407. Thompson, M. P. & Kurzrock, R. Epstein-Barr virus and cancer. Clin. Cancer Res. 44110, 803–21 (2004).

442 8. Mesri, E. A., Feitelson, M. A. & Munger, K. Human viral oncogenesis: a cancer hallmarks analysis. Cell Host Microbe 15, 266–82 (2014).

444 9. zur Hausen, H. Oncogenic DNA viruses. Oncogene 20, 7820–3 (2001).

445 10. PCAWG Consortium. Pan-cancer analysis of whole genomes. Nature (2019). 446 doi:10.1101/162784

447 11. Borozan, I. et al. CaPSID: A bioinformatics platform for computational pathogen 448 sequence identification in human genomes and transcriptomes. BMC Bioinformatics 13, 1–11 449 (2012).

450 12. Borozan, I., Watt, S. N. & Ferretti, V. Evaluation of alignment algorithms for 451 discovery and identification of pathogens using RNA-Seq. PLoS One 8, e76935 (2013).

452 13. Nicoll, M. P. et al. The HSV-1 Latency-Associated Transcript Functions to Repress 453 Latent Phase Lytic Gene Expression and Suppress Virus Reactivation from Latently Infected 454 Neurons. PLoS Pathog. 12, e1005539 (2016).

455 14. Krug, L. T. & Pellett, P. E. Roseolovirus molecular biology: recent advances. Curr. 456 Opin. Virol. 9, 170–7 (2014).

457 15. Newman, A. M. et al. Robust enumeration of cell subsets from tissue expression 458 profiles. Nat. Methods 12, 453–457 (2015).

459 16. Eliassen, E. et al. Human Herpesvirus 6 and Malignancy: A Review. Front. Oncol. 8, 460 512 (2018).

461 17. Spandole, S., Cimponeriu, D., Berca, L. M. & Mih escu, G. Human anelloviruses: an 462 update of molecular, epidemiological and clinical aspects. Arch. Virol. 160, 893–908 (2015).

463 18. van de Berg, P. J. et al. Human cytomegalovirus induces systemic immune activation 464 characterized by a type 1 cytokine signature. J. Infect. Dis. 202, 690–9 (2010).

465 19. Garcia-Martinez, A. et al. Lack of cytomegalovirus detection in human glioma. Virol. 466 J. 14, 216 (2017).

467 20. Fujimoto, A. et al. Whole-genome sequencing and comprehensive variant analysis of 468 a Japanese individual using massively parallel sequencing. Nat. Genet. 42, 931–6 (2010).

469 21. Furuta, M. et al. Characterization of HBV integration patterns and timing in liver 470 cancer and HBV-infected livers. Oncotarget 9, 25075–25088 (2018).

471 22. Canisius, S., Martens, J. W. M. & Wessels, L. F. A. A novel independence test for 472 somatic alterations in cancer shows that biology drives mutual exclusivity but chance 473 explains most co-occurrence. Genome Biol. 17, 261 (2016).

474 23. Kawai-Kitahata, F. et al. Comprehensive analyses of mutations and hepatitis B virus 475 integration in hepatocellular carcinoma with clinicopathological features. J. Gastroenterol. 476 51, 473–486 (2016).

477 24. Borozan, I., Zapatka, M., Frappier, L. & Ferretti, V. Analysis of Epstein-Barr Virus 478 Genomes and Expression Profiles in Gastric Adenocarcinoma. J. Virol. 92, e01239-17 479 (2018).

480 25. Mork, J. et al. Human Papillomavirus Infection as a Risk Factor for Squamous-Cell 481 Carcinoma of the Head and Neck. N. Engl. J. Med. 344, 1125–1131 (2001).

482 26. Li, N. et al. Human papillomavirus infection and bladder cancer risk: A meta-analysis. 483 J. Infect. Dis. 204, 217–223 (2011).

484 27. Cao, S. et al. Divergent viral presentation among human tumors and adjacent normal 485 tissues. Sci. Rep. 6, 28294 (2016).

486 28. Travé, G. & Zanier, K. HPV-mediated inactivation of tumor suppressor p53. Cell 487 Cycle 15, 2231–2 (2016).

488 29. Werness, B. A., Levine, A. J. & Howley, P. M. Association of human papillomavirus 489 types 16 and 18 E6 proteins with p53. Science 248, 76–9 (1990).

490 30. Scheffner, M., Werness, B. A., Huibregtse, J. M., Levine, A. J. & Howley, P. M. The 491 E6 oncoprotein encoded by human papillomavirus types 16 and 18 promotes the degradation 492 of p53. Cell 63, 1129–36 (1990).

493 31. Henderson, S., Chakravarthy, A., Su, X., Boshoff, C. & Fenton, T. R. APOBEC-494 Mediated Cytosine Deamination Links PIK3CA Helical Domain Mutations to Human 495 Papillomavirus-Driven Tumor Development. Cell Rep. 7, 1833–1841 (2014).

496 32. Burns, M. B., Temiz, N. A. & Harris, R. S. Evidence for APOBEC3B mutagenesis in 497 multiple human cancers. Nat. Genet. 45, 977–983 (2013).

498 33. Schlecht, N. et al. Gene expression profiles in HPV-infected head and neck cancer. J. 499 Pathol. 213, 283–293 (2007).

34. Nelson, P. N. et al. Demystified. Human endogenous retroviruses. Mol. Pathol. 56, 11–18 (2003).

35. Paces, J. et al. HERVd: the Human Endogenous RetroViruses Database: update. Nucleic Acids Res. 32, D50 (2004).

36. Pavlícek, A., Paces, J., Elleder, D. & Hejnar, J. Processed pseudogenes of human endogenous retroviruses generated by LINEs: their integration, stability, and distribution. Genome Res. 12, 391–9 (2002).

37. Bao, W., Kojima, K. K. & Kohany, O. Repbase Update, a database of repetitive elements in eukaryotic genomes. Mob. DNA 6, 11 (2015).

38. Ohnuki, M. et al. Dynamic regulation of human endogenous retroviruses mediates factor-induced reprogramming and differentiation potential. Proc. Natl. Acad. Sci. 111, 12426–12431 (2014).

39. Smith, C. C. et al. Endogenous retroviral signatures predict immunotherapy response in clear cell renal cell carcinoma. J. Clin. Invest. 128, 4804–4820 (2018).

40. Tang, K.-W. & Larsson, E. Tumour virology in the era of high-throughput genomics. Philos. Trans. R. Soc. Lond. B. Biol. Sci. 372, 20160265 (2017).

41. Jiang, Z. et al. The effects of hepatitis B virus integration into the genomes of hepatocellular carcinoma patients. Genome Res. 22, 593–601 (2012).

42. Hu, Z. et al. Genome-wide profiling of HPV integration in cervical cancer identifies clustered genomic hot spots and a potential microhomology-mediated integration mechanism. Nat. Genet. 47, 158–163 (2015).

43. Zhao, L.-H. et al. Genomic and oncogenic preference of HBV integration in hepatocellular carcinoma. Nat. Commun. 7, 12992 (2016).

44. Li, X. et al. The function of targeted host genes determines the oncogenicity of HBV integration in hepatocellular carcinoma. J. Hepatol. 60, 975–84 (2014).

45. Shen, C.-J., Cheng, Y.-M. & Wang, C.-L. LncRNA PVT1 epigenetically silences miR-195 and modulates EMT and chemoresistance in cervical cancer cells. J. Drug Target. 25, 637–644 (2017).

46. 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, 1–9 (2013).

47. Nault, J.-C. et al. Recurrent AAV2-related insertional mutagenesis in human hepatocellular carcinomas. Nat. Genet. 47, 1187–93 (2015).

48. Borozan, I. & Ferretti, V. CSSSCL: A python package that uses combined sequence similarity scores for accurate taxonomic classification of long and short sequence reads. Bioinformatics 32, 453–455 (2015).

49. Sung, W. K. et al. Genome-wide survey of recurrent HBV integration in hepatocellular carcinoma. Nat. Genet. 44, 765–769 (2012).

50. Alexandrov, L. B., Nik-Zainal, S., Wedge, D. C., Campbell, P. J. & Stratton, M. R. Deciphering Signatures of Mutational Processes Operative in Human Cancer. Cell Rep. 3, 246–259 (2013).

51. Wallace, N. A. & Münger, K. The curious case of APOBEC3 activation by cancer-associated human papillomaviruses. PLoS Pathog. 14, e1006717 (2018).

52. Roberts, S. A. et al. An APOBEC cytidine deaminase mutagenesis pattern is widespread in human cancers. Nat. Genet. 45, 970–976 (2013).

# Figure Legends

**Fig. 1: Overview, design and summary statistics.** (a) Workflow to identify and characterize viral sequences from the whole-genome and RNA sequencing of tumor and non-malignant samples. Viral hits were characterized in detail using several clinical annotations and resources generated by PCAWG. The red line represents the median. (b) Identified viral hits in contigs showing higher PMER's (viral reads **p**er **m**illion **e**xtracted **r**eads) for artificial sequences like vectors than the virus. Displayed are all viruses that occur in at least 20 primary tumor samples in the same contig together with an artificial sequence. (c) Summary of the viral search space used in the analysis grouped by virus genome type. The number of

555virus positive tumor samples are indicated in the outer rings (PMER log scale for WGS and 556RNA sequencing data) as detected by any of the pipelines. Taxonomic relations between the 557viruses are indicated by the phylogenetic tree. dsDNA: double stranded DNA virus, dsDNA-558RT: double-stranded DNA reverse transcriptase virus, ssDNA: single-stranded DNA virus, 559ssRNA-RT: single-stranded RNA reverse transcriptase virus, ssRNA: single-stranded RNA 560virus, dsRNA: double-stranded RNA virus. Fraction of hits in WGS and RNA sequencing 561data are depicted as stacked barplot.
562

**563Fig. 2: Detected viruses: Consensus for detected viruses in whole genome and 564transcriptome sequences.** Number of genus hits among tumor samples for the three 565independent pipelines and the consensus set defined by evidence from multiple pipelines. (a) 566based on whole genome sequencing, (b) and based on transcriptome sequencing. (c) Heatmap 567showing the total number of viruses detected across various cancer entities. The sequencing 568data used for detection is indicated among the total number of hits (WGS= blue, RNA-569seq=green). The fraction of virus positive samples is shown on top and the type of non-570malignant tissue used in the analysis is indicated if more than 15% of the analyzed samples 571are from a respective tissue type (solid tissue, lymph node, blood or adjacent to primary 572tumor). (d) t-SNE clustering of the tumor samples based on PMER of their consensus virome 573profiles, using Pearson correlation as the distance metric. Major clusters are highlighted by 574indicating the strongest viral genus and the dominant tissue types that are positive in that 575cluster. Dot size represents the viral reads per million extracted reads (PMER).
576

**577Fig. 3: Virus specific findings.** (a) HBV detections, validations and driver mutations in liver 578cancer. Star indicating mutual exclusivity between HBV detections and somatic driver gene 579mutations. Red boxes represent virus-positive tumor samples, purple - viral genomic 580integrations, green – driver mutations, grey – missing data. (b) Virus detections in gastric 581cancer samples, indication of virus phase (lytic/latent, dark red) and driver mutations (green). 582Yellow color indicates donors with virus-positive non-malignant samples. Grey box refers to 583samples with available RNA-seq data. (c) Virus detections (red) and driver mutations (green) 584in cervix (blue) and head and neck cancer (brown). Star indicating mutual exclusivity 585between alphapapillomavirus detections and somatic driver gene mutations. (d) 586Alphapapillomavirus detection and exposures of mutational APOBEC signatures SBS2 and 587SBS13, with sample sizes shown below. Wilcoxon rank-sum test (two-sided) revealed a 588significant difference ($P = 0.02$) of mutational signature exposure between virus-positive and 589negative head/neck tumor samples. Black line indicates median in each group. (e) Gene 590expression based tSNE map of head and neck cancer samples show a distinct gene expression 591profile for virus positive samples. Virus-positive and negative samples were labeled as red 592and grey dots, respectively. (f) The violin plot of APOBEC3B gene expression for 593alphapapillomavirus positive and negative samples in cervix and head/neck cancer (FDR 594corrected Wilcoxon rank-sum test, two-sided, $P = 1.6 \times 10^{-4}$). The center line represents 595median, the upper and lower boundaries of the violin plot refer to the maximum and 596minimum values, respectively. (g) Tumor-infiltrating immune cells as quantified by 597CIBERSORT using RNA-seq samples from head and neck cancer patients. All four cell types 598showed significant enrichment of immune cells in virus positive samples (FDR corrected 599Wilcoxon rank-sum test two-sided, n=24 vs 18). Tukey boxplot indicates the median by the 600middle line and the 25-75th percentiles by the box. The whiskers were drawn up to the 1.5 601interquartile range from the lower and upper quartile.
602

**603Fig. 4: Endogenous retroviruses**. (a) Heatmap showing the HERV expression across all 604tumor samples. HERV TPMs were grouped by family and summed up. Hierarchical 605clustering was performed by family based on Manhattan distance with complete linkage after 606log2 transformation of HERVs transcripts per million (TPM) expression values. (b) Fraction 607of active loci in the genome with a TPM >0.2 plotted against the fraction of samples. (c) 608TPM based expression of the highly expressed HERVs ERV1 and ERVK across tumor types. 609n described number of tumor samples analyzed. Violin plots marked with the median as red 610dot. The upper and lower boundaries of the violin plot extend out to the maximum and 611minimum values. (d) Survival difference between kidney cancer samples expressing high 612(red) and low levels (blue) of ERV1. Kaplan-Meier curve shows the overall survival of

613patients (n=113) with high and low levels of ERV1 using a cutoff of 16.3 tpm (Log-rank test 614$P$=0.0081). Patients at risk are provided below.

615

616**Fig. 5: Impact of virus integration**. (a) Integration sites detected in gene regions (including 617promoter, exon, intron and 5' UTR regions) are labeled in red for increased gene expression 618and blue for expression measured. Rows of each heatmap designate nearest genes to the 619integration sites and columns represent individual ICGC donor and project IDs. Intragenic 620HBV integration sites detected in liver cancers (ICGC project codes: LIRI, LIHC and LINC). 621For TERT and SEMA6D intergenic integrations are shown as well. (b) Integration sites 622detected for HPV-16 and 18 in head/neck (samples color coded magenta) and cervical 623(samples color coded blue) cancers (ICGC project codes: HNSC and CESC) gene labels with 624star indicated HPV18 as opposed to HPV16 viral integrations. (c) A local increase in the 625number of SCNAs was shown in the vicinity of HBV viral integrations (n=21 viral 626integrations in individual patients, $P$=7.4 × 10$^{-3}$; two-sided paired $t$-test). (d) Genomic 627visualization of the HBV virus integration sites relative to the *TERT* gene in five liver tumor 628patients. (e) The increased gene expression (FPKM) of *TERT* gene in two liver tumors with 629HBV viral integrations in comparison to the *TERT* expression in tumor and non-malignant 630adjacent tissue. Tumor samples with a non-coding driver mutation were labeled in orange.

631

# 632Methods

## 633Identifying potential pathogenic reads

634To reduce the number of reads to be considered for the pathogen search, we identified 635potential pathogenic reads using script available at https://github.com/mzapatka/p-dip. Based 636on the reads aligned by BWA[53] or STAR[54] to hg19 using the standard PCAWG approach, we 637identified read pairs where at least one read did not show a good mapping to the human 638genome (longest stretch of mapped bases from 20 to 30 bases), were unmapped or mapped to 639NC_007605 (human herpesvirus 4, which is contained in the 1000 genomes version of the 640hg19 human reference genome) and extracted these for further processing. To speed up the 641extraction, we used bamcollate2 from Biobambam2[55] v2.08 as input stream to the python 642script.

## 643Identification of endogenous retroviruses

644The expression of the endogenous retroviruses was analyzed based on the RNA sequencing 645data and aligned STAR based on the setting developed within PCAWG (hg19 and Gencode 64619). In contrast to the standard pipeline, the reference transcripts from Gencode 19 were 647enriched by adding HERV locations extracted from RepeatMasker (URL: 648http://www.repeatmasker.org, rmsk from UCSC, version 17/08/03) and Featurecounts 649(subread-1.5.3)[56] applied to identify reads mapping to the modified reference transcripts. 650Resulting reads counts were converted into transcripts per million (TPM) according to 651Wagner et al.[57].

## 652Norwich SEarching for PATHogens (SEPATH) pipeline

653Our starting point is to take reads that are not mapped to the human genome using the 654extracted potential pathogenic reads. Low quality bases (q<30) are trimmed from the read 655ends and the TruSeq indexed adapter and TruSeq universal adapter are removed using 656Cutadapt v1.8.1[58]. Reads less than 32 bp were discarded. Additional filtering is performed to 657remove reads containing more than 5% of Ns or those with low complexity (dust method 658with maximum score of 10) using Prinseq v0.20.3[59]. Metagenomic Phylogenic Analysis 659(MetaPhlAn)[60,61] is then applied to identify and quantify the presence of bacterial and viral 660populations. MetaPhlAn comes with a curated marker database of ~1M unique clade-specific 661marker genes identified from reference genomes (version 2.0 of the database was used).

662Reads are aligned against the unique marker gene database using BowTie2 v2.2.1[62] with
663presets set to sensitive. Reads are then counted and normalized giving a relative-abundance
664estimation at each level of the phylogenetic tree.

665Detection and Analysis of Microbial Infectious Agents by NGS P-DiP - 
666            Pathogen discovery pipeline

667The assembly based pipeline (P-DiP) was further developed based on a version implemented
668by Malik Alawi and Adam Grundhoff[63]. In summary, the pipeline runs preprocessing,
669assembly and BLAST searches and stores processing details and final results in a postgresql
670database. For the whole genome sequencing and RNA sequencing, we started with the
671potential pathogenic reads extracted from the BWA aligned whole genome sequencing bam
672files. As a first step, reads were trimmed based on quality using trimmomatic. Thereafter, host
673reads were subtracted by aligning to the human reference genome (WGS: hg19 excluding
674NC_007605    and    hs37d5    and    adding    phiX,    RNAseq:
675Homo_sapiens.GRCh37.dna.primary_assembly) using Bowtie2 (2.2.8)[62]. Trinity (v2.0.6)[64]
676was used for the read assembly of WGS reads which were not aligned by bowtie with
677sufficient quality (not aligned with --very-fast (-D 5 -R 1 -N 0 -L 22 -i S,0,2.50) to
678Homo_sapiens.GRCh37.ncrna, Homo_sapiens.GRCh37.cdna.all or PhiX) for the RNA
679sequencing data we applied idba assembler (V1.1.3)[65]. Assembled contigs were filtered by
680size (minimal length of 300 bp). Abundance was estimated by remapping all reads not
681aligning to the human reference to the assembled contigs using bowtie2 again. Putative PCR
682duplicates identified by mapping location were removed from the abundance count. The
683taxonomic classification of the size filtered contigs was performed using the BLAST+
684package (2.2.30)[66] and nucleotide databases nt (2015-05-15) and nr (2015-04-20). For the
685extraction of pathogen hits R-scripts were used to filter the blast results (at
686https://github.com/mzapatka/p-dip). In summary, for each of the contig, the best BLAST hits
687for each segment of the contig were considered and the reads aligning to these segments
688identified. Potential contaminants were defined based on the taxonomy annotation in NCBI
689taxonomy. Any taxonomy id below plasmids (36549), transposons (2387), midivariant
690sequences (31896), insertion sequences (2673), artificial sequences (81077) and synthetic
691viruses (512285) was annotated as potential contamination. Segments with higher read counts
692of these sequences compared to pathogen hits were flagged as contaminants and not further
693considered.

694Computational Pathogen Sequence Identification (CaPSID) description of the
695            analysis workflow

696CaPSID's[11] metagenomic analysis pipeline starts by first processing a BAM file containing
697reads sequenced from a tumor (or normal) sample aligned to the human reference sequence
698(GRCh37/hg19). Reads that did not map to the human reference are extracted and filtered for
699low complexity and quality using the SGA[67] preprocessing module and then aligned in single-
700end mode using the Bowtie2 aligner[62] to 5,652 NCBI[68] viral reference sequences (RefSeq)
701and a filter sequence reference database composed of 5,242 bacterial and 1,138 fungal
702reference sequences also downloaded from the NCBI. In order to improve the sensitivity and
703specificity with which viral sequences are detected, reads that did not map to any reference
704with Bowtie2 are realigned against the same viral RefSeq database, using a more sensitive
705SHRiMP2 aligner using its local alignment mode[69]. At the completion of this two-step
706alignment process, reads aligning to viral reference sequences are annotated using the
707information stored in the CaPSID's genome database containing full NCBI GenBank and taxa
708information. Using information from each aligned read CaPSID then calculates the following
709four metrics: (i) the total number of reads (or hits) aligning across any given viral genome,
710(ii) the total number of reads aligning only across gene regions within any given viral
711genome, (iii) the total coverage across each viral genome and (iv) the maximum coverage
712across any of the genes in a given viral genome.
713
714*Filtering of viral candidates with low significance*

715

716 In a typical analysis of tumor whole genome or transcriptome sample, CaPSID reports
717 candidate sequences from dozens of different viral genomes, some of which are not related to
718 cancer phenotype. Some of these reported viral hits are also due to a series of experimental
719 and computational artifacts. In order to reduce the number of potential false positives CaPSID
720 pipeline flags viral genomes could be the result of artifacts present in the sequencing data or
721 those with no obvious relation to cancer phenotype and that could be filtered later on. The
722 following criteria are used to flag and filter for potential viral candidates: (i) flag viral
723 candidates with low coverage, (ii) flag bacteriophage viral genome sequences, (iii) report
724 only viral candidates with read composition different from the one expected when generated
725 from the host's reference GRCh37/hg19 sequence, (iv) flag viral candidates that are typically
726 not known to infect humans and those with low read abundance and/or low overall alignment
727 read accuracy.

728

729 In the first step CaPSID flags viral genomes with low read count and/or coverage using its
730 three metrics including: total number of uniquely aligned reads < 3, total genome coverage <
731 10% and maximum gene coverage < 50%. Viral genomes with low read count can arise as a
732 result of i) low read/transcript abundance in the human sequenced sample, ii) non-specific
733 alignment between sequenced short reads (for example low complexity reads) and viral
734 reference sequences and iii) for RNA-seq library preparation where highly expressed
735 transcripts generally dominate over low abundance targets. In order to limit reporting viral
736 genomes with very low coverage, we chose to flag all those with maximum gene coverage <
737 50%. Since this lower bound on the maximum gene coverage applies to individual genes and
738 not to the complete viral genome, it is unlikely that viruses with such low coverage are
739 biologically significant. The second step in our filtering approach is to flag bacteriophage
740 viral genomes that are most likely not related to any cancer phenotype. Bacteriophages are
741 detected as a result of the presence of bacteria (or bacterial contamination) in human
742 sequenced samples. The third step is used to determine whether the genome coverage
743 observed for each viral candidate is different from the one expected to arise from reads
744 originating exclusively from the human reference DNA GRCh37/hg19 sequence. To build the
745 CaPSID background model we use the ART NGS read simulator. The entire GRCh37/hg19
746 sequence reference file is first fed to the ART[70] simulator (parameters: art_illumina [Illumina
747 platform] -l [read length = 100 bp] -f [the fold of read coverage to be simulated = 100] with
748 default values for indels and substitution rates), which then generates single-end (or paired-
749 end) reads and base quality values.

750

751 Reads simulated by ART are then aligned to the viral reference sequence database using the
752 same alignment approach for reads originating from tumor samples (see above). CaPSID then
753 calculated the four metrics for the GRCh37/hg19 background model using the alignment
754 information from simulated reads aligning to viral reference sequences. The fourth step
755 consists of flagging viral candidates that are typically not known to infect humans using a
756 dictionary of ~ 130 terms that we have compiled from a database of all viruses known to
757 infect humans. In addition to the above filtering criteria CaPSID also considers the read
758 abundance associated with each viral candidate sequence (abundance is expressed in terms of
759 aligned reads in parts-per million of total number of unmapped reads) and the average read
760 percent identity with which reads align to a given viral candidate reference sequence.

761

762 *De novo assembly and taxonomic classification of contigs*

763

764 The purpose of this analysis step is to attempt to characterize potential novel viral sequences
765 at the species or subspecies level. Unaligned reads which could not be aligned to any of the
766 filter/host or viral reference sequences are assembled into contigs using the IDBA algorithm[65].
767 Assembled contigs are then masked for repeat regions using RepeatMasker and then filtered
768 for their size and read coverage (contig length >= 500 bp and coverage > 5x). Resulting
769 contigs are then assigned into taxonomic groups at the genus level using the CSSSCL
770 algorithm[48]. Contigs lacking sequence homology to reference sequences contained in the
771 CaPSID or blast nucleotide databases with percent identity < 90% are then selected as
772 suggestive of the presence of new viral strains/isolates or species.

## 773Defining consensus hits

774Identification of the consensus hits was achieved by optimizing two features of the individual 775genus hits: PMER 1 as cutoff (see analysis of the validation set) and percentage identity 776>90%. 90% percentage identity threshold was determined based on our benchmarking study[12] 777indicating that an alignment-based approach can still accurately characterize viral sequences 778with up to 10% mutation rate (when compared to sequences stored in a reference database). 779Lowering the threshold, with which short reads align to any given reference sequence below 78090% identity on average, results in a drop of sequence coverage due to a high attrition rate of 781aligned reads, lowering the detection rate and thus providing more uncertain characterization 782of viral candidates. Notably, there was no difference in the PMER distribution of common 783hits across the three pipelines indicating that a common detection cut-off is reasonable 784(Extended Data Figure 3b).
785
786The consensus set was restricted to genera that were covered in at least two detection 787pipelines (Extended Data Figure 1b). Notably, we could not detect any more hits with high 788PMERs using the unique search space of P-DiP, indicating that almost all of the viral hits 789from individual pipelines were also screened by another pipeline.


## 790Virus integration detection analysis

791A subset of viral candidates identified to be present in tumor samples by the CaPSID's 792analysis pipeline (parameters used: PMER >= 1.1 and genome coverage > simulated 793background model) was selected for the detection of viral integration events using the 794VERSE[71] algorithm. This subset of viruses included: Herpesviruses (HHV-1, 2, 4, 5, 6A/B), 795Simian virus 40 (SV40) and 12 (SV12), Human immunodeficiency virus (HIV1), Human and 796Simian T-cell lymphotropic virus type 1 (HTLV1 and STLV1), BK polyomavirus (BKP), 797Human parvovirus B19, Mouse mammary tumor virus, Murine type C retrovirus, Mason-798Pfizer monkey virus, Hepatitis B (HBV), Papilloma viruses (HPV-16, 18 and 6a and Adeno-799associated virus - 2 (AAV2). Below we describe the steps used for viral integration detection 800analysis.
801
802Viral integration events in the host can be detected using paired-end NGS technologies that 803facilitate the detection of genomic rearrangements, as well as gene fusions and novel 804transcripts. VERSE is capable of determining virus integration sites within a single base 805resolution by requiring the presence of both chimeric and soft clipped reads. In addition, 806VERSE improves the detection through customizing reference genomes and was shown to 807substantially enhance the sensitivity of virus integration site detection[71]. VERSE categorizes 808its predictions into one of two classes: (a) a 'high' confidence hit with a single base resolution 809- if there is a sufficient number of soft-clipped reads to support an integration locus so that 810CREST is able to detect it; (b) a 'low' confidence hit with a 10 bp resolution where CREST 811has failed to detect an integration event for the lack of high quality soft-clipped reads.
812
813In order to further limit the false positive rate associated with viral integration sites we 814compare results obtained with VERSE to those from Fujimoto et al[72]. Out of 64 whole 815genome liver cancer samples with HBV integration events reported in Fujimoto et al., 50 are 816part of the PCAWG dataset analyzed in this study. 45 out of 50 of these samples tested 817positive for HBV when analyzed by CaPSID (filtering criteria used; PMER >= 1.0, genome 818coverage > host background model and read % identity >= 89%). In addition, 50 of these 819WGS samples had 23 matching whole transcriptome (WT) samples and 22 of these were 820identified to be positive for HBV by CaPSID (filtering criteria used; maximum gene coverage 821>= 50%, read % identity >= 89% and PMER >= 1.0). By combining WGS and whole 822transcriptome tumor samples together, 47/50 in total tested positive for HBV when analyzed 823by CaPSID.
824
825Using VERSE, virus integration sites were detected in 28/47 (60%) of these. This result 826indicates that for a subset of viral integration events, VERSE might be a more stringent 827approach when compared to the one used in Fujimoto et al. This can be explained by the fact

828that VERSE requires both the presence of paired-end chimeric and soft clipped reads while 829the method presented in Fujimoto et al. relies on paired-end reads only. In order to explore 830these results further we compared integration sites obtained with VERSE and Fujimoto et al. 831with an overlapping window of 10 bp. Our analysis indicates that among 23 integration sites 832identified by VERSE in whole transcriptome data and that overlap with the results from 833Fujimoto et al., 91% of these are classified with high confidence hits and only 9% with low 834(N total overlap = 23, high = 21 (91%) and low = 2 (9%)). However, a similar result is not 835observed for integration events found using WGS data (N total overlap = 14, high = 6 (43%), 836low = 8 (57%)) where the proportion of integration events classified as high and low is 837similar.

838Thus, our analysis indicates that one important factor for improving the agreement between 839these two datasets is the confidence level assigned by VERSE to each candidate integration 840site - but only in the case when integration sites are detected using whole transcriptome data. 841In order to reduce the potential number of false positives we decided to use all integration 842sites predicted by VERSE when these are obtained using WGS data and only high confidence 843calls when using whole transcriptome data.

## 844Contaminations

845Based on the presence of vector sequences in the contig assembled by the P-DiP and based on 846the background model from CaPSID we could identify which virus hits originate from 847common lab contaminants or due to sequence similarities to the human genome. In addition, 848we filtered known contaminants (see below). For P-DiP we filtered all hits not having more 849target reads than any artificial sequence (excluding artificial viruses) on an individual contig 850region. Hits caused by vector and other artificial sequences were identified analyzing the 851assembled contigs for combined hits to viral pathogens and artificial sequences. Checking 852viral hits occurring at least 40 times in a such contig we could clearly separate contaminants 853from viral pathogens.

854The gammaretrovirus hits (NCBI taxonomy id: 153135, species: murine leukemia virus) were 855also marked as artifacts, based on the additional BLAST hits of the corresponding contigs to 856the *Mus musculus* genome by P-DiP, as well as on the background model of the CaPSID 857pipeline designed to limit the number of spurious hits. Most frequent virus hits prone to 858contamination by artificial sequences are Lambdalikevirus, Alphabaculovirus, Microvirus, 859Simplexvirus, Hepacivirus, Cytomegalovirus, Orthopoxvirus and Punalikevirus. But 860restricting to at least 1 PMER for the potential virus hit contaminants drop to one 861Cytomegalovirus case.

## 862Filtering contaminants

863We filtered all Microviridae (taxonomy ID: 10841) because of the phix174 spike-in used 864during sequencing. Caudovirales (taxonomy ID: 28883), tailed bacteriophages, were removed 865as they typically infect bacterial hosts. Baculoviridae were filtered because of infecting insect 866cells and commonly being used in the lab. The virus coverage was analyzed by aligning the 867potential pathogenic reads with BWA mem to the human hg19 reference genome after adding 868the respective virus reference sequence most frequently detected within the genus. Coverage 869was thereafter calculated base specific using BEDTools coverage. As we identified EBV in 870all 14 normal blood controls from ovarian cancer that were EBV immortalized these were 871removed from the virus hits.

## 872Integration of external PCAWG datasets

873We tested for mutual exclusivity e.g between virus detections and driver gene mutations by 874applying DISCOVER[22]. Based on the gene expression data, immune cell proportions were 875analyzed by CIBERSORT[15]. For survival analysis, Cox proportional hazards analysis was 876performed using R libraries 'survival' and 'survminer' for the figures. The optimal cutpoints 877were identified by maxstat using the method presented in Lausen and Schumacher[73] (library 878maxstat).
879

## 880 Virus load

881 The viral load in relation to the human genome equivalents was calculated based on the 882 human bases sequenced (read length x number of reads mapped to the human genomes), 883 tumor sample purity (if available 100% otherwise) assuming a ploidy of two and using a 884 human genome size of 2,897,310,462 bases (mappable part of the human genome). This 885 number of human genome equivalents was then related to the viral genome equivalents 886 calculated based on viral reads identified, read length and virus genome size.

887 $$tumor\ genome\ equivalents = \frac{read\ length \times number\ of\ reads\ mapped\ to\ the\ human\ genomes}{mappable\ human\ genome\ size\ \times\ tumor\ ploidy} \times tumor\ purity$$

888 $$virus\ genome\ equivalents = \frac{read\ length \times number\ of\ viral}{virus\ genome\ size}$$

889 $$virus\ load = \frac{virus\ genome\ equivalents}{tumor\ genome\ equivalents}$$

## 890 Human research participants

891 The Ethics oversight for the PCAWG protocol was undertaken by the TCGA Program Office 892 and the Ethics and Governance Committee of the ICGC. Each individual ICGC and TCGA 893 project that contributed data to PCAWG had their own local arrangements for ethics 894 oversight and regulatory alignment.

## 895 Statistics

896 If not specified otherwise, we used two-sided Wilcoxon rank-sum test for groups with n >3.
897 Further details can be accessed at the ' Life Sciences Reporting Summary'.
898

# 899 Data Availability Statement

900 Somatic and germline variant calls, mutational signatures, subclonal reconstructions, 901 transcript abundance, splice calls and other core data generated by the ICGC/TCGA Pan-902 cancer Analysis of Whole Genomes Consortium is described here[10] and available for 903 download at https://dcc.icgc.org/releases/PCAWG. Additional information on accessing the 904 data, including raw read files, can be found at https://docs.icgc.org/pcawg/data/. In 905 accordance with the data access policies of the ICGC and TCGA projects, most molecular, 906 clinical and specimen data are in an open tier which does not require access approval. To 907 access potentially identification information, such as germline alleles and underlying 908 sequencing data, researchers will need to apply to the TCGA Data Access Committee (DAC) 909 via dbGaP (https://dbgap.ncbi.nlm.nih.gov/aa/wga.cgi?page=login) for access to the TCGA 910 portion of the dataset, and to the ICGC Data Access Compliance Office (DACO; 911 http://icgc.org/daco) for the ICGC portion. In addition, to access somatic single nucleotide 912 variants derived from TCGA donors, researchers will also need to obtain dbGaP 913 authorization.
914 Data sets described specifically in this manuscript can be found in the Supplementary Tables.
915

# 916 Code availability Statement

917 The core computational pipelines used by the PCAWG Consortium for alignment, quality 918 control and variant calling are available to the public at https://dockstore.org/search? 919 search=pcawg under the GNU General Public License v3.0, which allows for reuse and 920 distribution. The pathogen discovery pipeline P-DiP is available on github at 921 https://github.com/mzapatka/p-dip. CaPSID is available from the github pages (
922 pipeline: https://github.com/capsid/capsid-pipeline,
923 webapp: https://github.com/capsid/capsid-webapp). The taxonomic classifier CSSSCL is 924 available from https://github.com/oicr-ibc/cssscl.

925

# 926 Methods-only References

927 53. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-928 MEM. arXiv Prepr. arXiv 1303.3997 (2013).

929 54. Dobin, A. et al. STAR: ultrafast universal RNA-seq aligner. Bioinformatics 29, 15–21 930 (2013).

931 55. Tischler, G. & Leonard, S. Biobambam: Tools for read pair collation based algorithms 932 on BAM files. Source Code Biol. Med. 9, 1–17 (2014).

933 56. Liao, Y., Smyth, G. K. & Shi, W. FeatureCounts: An efficient general purpose 934 program for assigning sequence reads to genomic features. Bioinformatics 30, 923–930 935 (2014).

936 57. Wagner, G. P., Kin, K. & Lynch, V. J. Measurement of mRNA abundance using 937 RNA-seq data: RPKM measure is inconsistent among samples. Theory Biosci. 131, 281–285 938 (2012).

939 58. Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing 940 reads. EMBnet.journal 17, 10 (2011).

941 59. Schmieder, R. & Edwards, R. Quality control and preprocessing of metagenomic 942 datasets. Bioinformatics 27, 863–4 (2011).

943 60. Truong, D. T. et al. MetaPhlAn2 for enhanced metagenomic taxonomic profiling. Nat. 944 Methods 12, 902–3 (2015).

945 61. Segata, N. et al. Metagenomic microbial community profiling using unique clade-946 specific marker genes. Nat. Methods 9, 811–4 (2012).

947 62. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. Nat. 948 Methods 9, 357–9 (2012).

949 63. Fischer, N. et al. Rapid Metagenomic Diagnostics for Suspected Outbreak of Severe 950 Pneumonia. Emerg. Infect. Dis. 20, 1072–1075 (2014).

951 64. Grabherr, M. G. et al. Full-length transcriptome assembly from RNA-Seq data 952 without a reference genome. Nat. Biotechnol. 29, 644–652 (2011).

953 65. Peng, Y., Leung, H. C. M., Yiu, S. M. & Chin, F. Y. L. IDBA-UD: a de novo 954 assembler for single-cell and metagenomic sequencing data with highly uneven depth. 955 Bioinformatics 28, 1420–8 (2012).

956 66. Camacho, C. et al. BLAST+: architecture and applications. BMC Bioinformatics 10, 957 421 (2009).

958 67. Simpson, J. T. & Durbin, R. Efficient de novo assembly of large genomes using 959 compressed data structures. Genome Res. 22, 549–56 (2012).

960 68. Pruitt, K. D., Tatusova, T., Klimke, W. & Maglott, D. R. NCBI Reference Sequences: 961 current status, policy and new initiatives. Nucleic Acids Res. 37, D32–D36 (2009).

962 69. David, M., Dzamba, M., Lister, D., Ilie, L. & Brudno, M. SHRiMP2: Sensitive yet 963 Practical Short Read Mapping. Bioinformatics 27, 1011–1012 (2011).

964 70. Huang, W., Li, L., Myers, J. R. & Marth, G. T. ART: A next-generation sequencing 965 read simulator. Bioinformatics 28, 593–594 (2012).

966 71. Wang, Q., Jia, P. & Zhao, Z. VERSE: a novel approach to detect virus integration in 967 host genomes through reference genome customization. Genome Med. 7, 2 (2015).

968 72. Fujimoto, A. et al. Whole-genome mutational landscape and characterization of 969 noncoding and structural mutations in liver cancer. Nat. Genet. 48, 500–9 (2016).

970 73. Lausen, B. & Schumacher, M. Maximally Selected Rank Statistics. Biometrics 48, 971 73–85 (1992).

972

# 973 PCAWG-Pathogens Members

974 Malik Alawi[5,6], Ivan Borozan[2], Daniel S Brewer[3,4], Colin S Cooper[4,13], Nikita Desai[7,8], Roland Eils[14,15,16], 975 Vincent Ferretti[17,18], Adam Grundhoff[5], Murat Iskar[1], Kortine Kleinheinz[20,21], **Peter Lichter#**[1,10],

976 Hidewaki Nakagawa[22], Akinyemi I Ojesina[23,24,25], Chandra Sekhar Pedamallu[26,27,28], Matthias
977 Schlesner[20,29], Xiaoping Su[30] and Marc Zapatka[1]
978
979 # Corresponding author
980
981 20. Division of Theoretical Bioinformatics, German Cancer Research Center (DKFZ), Heidelberg 69120,
982 Germany.
983 21. Institute of Pharmacy and Molecular Biotechnology and BioQuant, Heidelberg University, Heidelberg
984 69120, Germany.
985 22. RIKEN Center for Integrative Medical Sciences, Yokohama, Kanagawa 230-0045, Japan.
986 23. Department of Epidemiology, University of Alabama at Birmingham, Birmingham, AL 35294, USA.
987 24. HudsonAlpha Institute for Biotechnology, Huntsville, AL 35806, USA.
988 25. O'Neal Comprehensive Cancer Center, University of Alabama at Birmingham, Birmingham, AL 35294,
989 USA.
990 26. Broad Institute of MIT and Harvard, Cambridge, MA 02142, USA.
991 27. Harvard Medical School, Boston, MA 02115, USA.
992 28. Department of Medical Oncology, Dana-Farber Cancer Institute, Boston, MA 02115, USA.
993 29. Bioinformatics and Omics Data Analytics, German Cancer Research Center (DKFZ), Heidelberg 69120,
994 Germany.
995 30. University of Texas MD Anderson Cancer Center, Houston, TX 77030, USA.
996

**a**

**PCAWG Input**

WGS: 5354 tumor/normal
35 cancer types

RNA-seq: 1057 tumors
25 cancer types

Number of analysed reads (million)

Thyroid
Lung
Biliary
Ovary
Kidney
Esophagus
Prostate
Head/Neck
Colon/Rectum
Pancreas
Bladder
Myeloid
Lymphoid
Cervix
Uterus
Bone/SoftTissue
Breast
Skin
CNS
Stomach
Liver

**Virus detection**

i. P-DiP
- Assembly based contig generation
- BLAST against nt and nr

ii. CaPSID
- Bowtie2/SHRiMP2 alignment to Genbank
- De novo assembly with unmapped reads, CSSSCL classifier for taxonomy assignment

iii. SEPATH
- Alignment to unique clade-specific marker genes for taxonomy assignment
- Relative-abundance estimation

**Virus integration sites**

i. VERSE
- Reference genome customization

**Integrative analysis with clinical data**

Consensus calls
2 out of 3 methods
PMER >1
356 positive donors
23 virus genera

APOBEC signature

Gene expression profiles

Impact on survival

Mutual exclusivity with cancer drivers

**Functional effects of virus integration**

Structural variants
SNVs
Expression changes

**b** Potential contaminants

**c**

dsDNA
dsDNA−RT
ssDNA
ssRNA−RT
ssRNA(−)
dsRNA

Viruses

WGS
286
RNA-seq
286

Fraction of hits (%)

WGS
RNA-seq

organ system:
Biliary
Bone/Soft Tissue
Breast
CNS
Cervix
Colon/Rectum
Esophagus
Head/Neck
Kidney
Liver
Lymphoid
Myeloid
Ovary
Pancreas
Prostate
Skin
Stomach
Thyroid
Uterus
Lung

**a** WGS (tumor) Number of hits | **b** RNA-seq (tumor) Number of hits

Legend:
- Evidence from multiple pipelines (red)
- Unique to one pipeline (gray)

**c**

Tumor

| Virus | Liver | Pancreas | Stomach | Head/neck | Esophagus | Cervix | Lymphoid | CNS | Breast | Colon/Rectum | Kidney | Prostate | Skin | Lung | Uterus | Ovary | Thyroid | Bladder | Bone/Soft | Biliary | Myeloid |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Lymphocryptovirus** (Epstein-Barr virus) | 13 | 10 | 26 | 16 | 23 | 2 | 12 | 4 | 2 | 10 | 6 | 5 | 7 | 3 | 2 | 2 | 1 | | | | 1 |
| **Orthohepadnavirus** (Hepatitis B virus) | 62 | | | | | | | 1 | 2 | | 1 | | | | | | | | | 1 | |
| **Roseolovirus** (Human herpesvirus 6B) | 2 | 21 | 6 | 2 | 3 | 1 | 2 | 4 | 3 | 5 | 4 | 2 | 1 | | 1 | 1 | | | | | |
| **Alphapapillomavirus** (Alphapapillomavirus 7 & 9) | | | | 18 | | 19 | | | | | | | | | | | | 2 | | | |
| **Cytomegalovirus** (Human cytomegalovirus) | | 2 | 13 | | 3 | | | 1 | | | | | | | | | | | | | |
| **Gammaretrovirus** (Murine leukemia virus) | 1 | 8 | 2 | | | | | | | | 1 | 3 | | | | 1 | | | | | |
| **Alphatorquevirus** (Torque teno virus 5) | | 1 | 2 | | | | 1 | 3 | | | | 1 | | | | | | | | | |
| **Lentivirus** (Human immunodeficiency virus) | | | | | | | 1 | 1 | | | 1 | | | 3 | | | | | 2 | | |

Non-malignant

| Virus | Liver | Pancreas | Stomach | Head/neck | Esophagus | Cervix | Lymphoid | CNS | Breast | Colon/Rectum | Kidney | Prostate | Skin | Lung | Uterus | Ovary | Thyroid | Bladder | Bone/Soft | Biliary | Myeloid |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Lymphocryptovirus | 6 | 1 | 19 | 5 | 2 | 1 | 2 | | 4 | 4 | 5 | 2 | 2 | 4 | 1 | 5 | 1 | | | | 3 |
| Orthohepadnavirus | | | | 1 | | | 1 | | | | | | | | | | | | | 1 | |
| Roseolovirus | 13 | 22 | 14 | 1 | 1 | 2 | | 2 | 5 | 3 | 15 | 4 | 1 | 1 | | 1 | | | | 3 | 2 |
| Alphapapillomavirus | | | | 1 | | | | | 1 | | | | | | | | | | | | |
| Cytomegalovirus | | | 11 | | | | | | | | | | | | | | | | | | 1 |
| Gammaretrovirus | | 3 | 1 | | | | 2 | | 1 | | 3 | | | | | | | | | | 1 |
| Alphatorquevirus | 13 | 1 | 1 | | | | 11 | 4 | 1 | 1 | 1 | 2 | | 1 | | | | | | | |
| Lentivirus | | | 2 | | | | | | | | | | | 2 | | 1 | | | | | |

Legend: WGS (blue) · RNA-seq (green) · WGS&RNA-seq (dark green) · Non-malignant: Solid tissue (red) · Blood derived (yellow) · Lymph node (green) · Tissue adjacent to primary (light blue)

**d**

Alphapapillomavirus
Cervix, head/neck, bladder

Roseolovirus
Pancreas

Gammaretrovirus
Pancreas

Lymphocryptovirus
Stomach, esophagus, colon

Orthohepadnavirus
Liver

Cytomegalovirus
Stomach

PMER: 1 10 20

tSNE1 / tSNE2

**a** Orthohepadnavirus-positive (HBV) samples in liver cancers (LIRI-JP, n=251)

Orthohepadnavirus (WGS+RNA) — n=50
PCR validation — n=50
Orthohepadnavirus (RNA) — n=18
HBAg — n=68
*TP53* — n=75
*CTNNB1* — n=56
*ARID1A* — n=47

Virus positive | Human genome integrations | Driver mutations | Not available

**b** Lymphocryptovirus-positive (EBV) samples in gastric cancers (n=75)

Lymphocryptovirus ( Lytic) — n=33
Cytomegalovirus — n=17
Roseolovirus — n=18
*TP53* — n=31
*ARID1A* — n=21
*CDKN2A* — n=11
*PIK3CA* — n=11

Virus positive | Driver mutations | RNA-seq available
Virus positive only in non-malignant solid tissue/adj. to tumor

**c** Cervix ● Head/Neck ●

Alphapapillomavirus — n=37
Lymphocryptovirus — n=18
*TP53* — n=30
*CDKN2A* — n=23
*TERT* — n=20
*FAT1* — n=17

Virus positive | Driver mutations

**d**
Mutational exposure of SBS2 [%]
$P = 0.02$

Mutational exposure of SBS13 [%]

n= 39 | 21 | 39 | 212 | 12 | 1
Alphapapillomavirus positive | Bladder | Head/Neck | Breast | Biliary | Cervix
Alphapapillomavirus negative

**e** Head/Neck gene expression profile
tSNE2 / tSNE1
n=42
Alphapapillomavirus
positive | negative

**f** *APOBEC3B*
gene expression (log2 FPKM)
n= 25 Cervix Head/Neck | 15 Cervix | 18 Head/Neck

**g** Head/Neck
CIBERSORT immune cell abundance
$P =$ 0.018 | 0.004 | 0.012 | 0.012
n= 24/18 | 24/18 | 24/18 | 24/18
Macrophages M1 | T cells CD8 | T cells follicular helper | T cells regulatory

**a**

tpm (log2)

−4 −2 0 2 4 6

Specimen
Histology

ERVL
ERVL-MaLR
ERV1
ERVK
ERV

**Specimen**
- Tumor
- Non-malignant
  (tissue adjacent to primary)

**Histology**
- Adenocarcinoma
- Chronic lymphocytic leukemia
- Diffuse glioma
- Hepatocellular carcinoma
- Liposarcoma, soft tissue
- Lobular carcinoma
- Mature B−cell lymphoma
- Melanoma
- RCC (distal tubules)
- RCC (proximal tubules)
- Sarcoma, soft tissue
- Squamous cell carcinoma
- Transitional cell carcinoma

**b**

HERV Family
- ERVL (n=102)
- ERVL−MaLR (n=80)
- ERV1 (n=260)
- ERVK (n=37)
- ERV (n=3)

Fraction of action of active sites (tpm>0.2) [%]

Samples [%]

**c**

ERV1

Expression [tpm]

ERVK

Expression [tpm]

Adenocarcinoma (n=341)
Chronic lymphocytic leukemia (n=31)
Diffuse glioma (n=46)
Hepatocellular carcinoma (n=49)
Liposarcoma, soft tissue (n=18)
Lobular carcinoma (n=6)
Mature B−cell lymphoma (n=101)
Melanoma (n=33)
RCC (distal tubules) (n=45)
RCC (proximal tubules) (n=69)
Sarcoma, soft tissue (n=15)
Squamous cell carcinoma (n=107)
Transitional cell carcinoma (n=23)

**d**

Survival probability

$P = 0.0081$

ERV1 Expression in Kidney (tpm > 16.3)
- ERV1-high
- ERV1-low

| | | | | |
|---|---|---|---|---|
| ERV1-high | 71 | 19 | 8 | 1 | 0 |
| ERV1-low | 42 | 30 | 20 | 7 | 1 |

Time (days)

**a** HBV integrations in ■ Liver cancer

**b** HPV16/18* integrations in ■ Cervix, ■ Head/Neck

Gene list (a): *TERT*, *KMT2B*, *SEMA6D*, *TEKT3*, *CCNA2*, *THRB*, *CDK15*, *BTD*, *CDH13*, *LIPI*, *MARS2*, *NCOR2*, *PUM1*, *DOK5*, *CPA6*, *CMIP*, *HEATR6*, *SFMBT2*, *RB1CC1*, *FLJ36000*, *C1orf198*, *ERICH1*, *GIMAP5*, *AARS2*, *CCNE1*, *SLC35F3*, *LINC01158*, *RGS12*, *FSIP2*, *USP9Y*, *PASD1*, *SLC1A7*, *ADAM3A*, *GRXCR1*, *MARK1*, *DKK2*, *ZMYM4*, *SENP5*

Gene list (b): *STX17*, *TEX10*, *DOLPP1*, *NR4A2**, *LINC00111*, *ETS2*, *PHLDB2**, *PLGRKT*, *ERBB2*, *PVT1*, *CEACAM5*, *MAMLD1*, *RASA3*, *FRMPD4*, *MAGI2*, *SLC9A7*, *CRAT*, *IQGAP1*, *IFT140*, *ENTPD1*, *SORBS1*, *COL6A6*, *STX17–AS1*, *ABLIM1**, *TALDO1**, *P3H2*, *UTP11*, *CASZ1**, *LINC–PINT**, *TP63*

■ Increased gene expression   Expression data ◨ available   ■ not available

**c** *P* = 7.4e-3

Number of SCNAs within 1Mbp

genomic background — virus integrations

**d** HBV virus integrations

*TERT*

chr5: 1,250,000 - 1,300,000

1,250 kb — 5 kb — 1,300 kb

**e** *TERT* FPKM (UQ normalized)

Liver gene expression profiles

tumor (n=105)

non-malignant (n=63)

▼ Virus integration   ● Non-coding driver mutation