

RESEARCH

Open Access



An atlas of the tissue and blood metagenome in cancer reveals novel links between bacteria, viruses and cancer

Sven Borchmann^{1,2,3}

Abstract

Background: Host tissue infections by bacteria and viruses can cause cancer. Known viral carcinogenic mechanisms are disruption of the host genome via genomic integration and expression of oncogenic viral proteins. An important bacterial carcinogenic mechanism is chronic inflammation. Massively parallel sequencing now routinely generates datasets large enough to contain detectable traces of bacterial and viral nucleic acids of taxa that colonize the examined tissue or are integrated into the host genome. However, this hidden resource has not been comprehensively studied in large patient cohorts.

Methods: In the present study, 3025 whole genome sequencing datasets and, where available, corresponding RNA-seq datasets are leveraged to gain insight into novel links between viruses, bacteria, and cancer. Datasets were obtained from multiple International Cancer Genome Consortium studies, with additional controls added from the 1000 genome project. A customized pipeline based on KRAKEN was developed and validated to identify bacterial and viral sequences in the datasets. Raw results were stringently filtered to reduce false positives and remove likely contaminants.

Results: The resulting map confirms known links and expands current knowledge by identifying novel associations. Moreover, the detection of certain bacteria or viruses is associated with profound differences in patient and tumor phenotypes, such as patient age, tumor stage, survival, and somatic mutations in cancer genes or gene expression profiles.

Conclusions: Overall, these results provide a detailed, unprecedented map of links between viruses, bacteria, and cancer that can serve as a reference for future studies and further experimental validation.

Introduction

Bacterial [1, 2] and viral [3–5] infections have widely been recognized as causes of cancer. Examples of carcinogenic viruses are Human Papillomaviridae, causing head and neck [6, 7] as well as cervical cancer [3, 8, 9] or *Hepatitis B virus*, causing liver cancer [10, 11]. The

main carcinogenic mechanisms for viral carcinogenesis are thought to be (1) viral integration into and disruption of the host genome and (2) expression of oncogenic viral proteins [12].

An important example of bacteria causing cancer is *Helicobacter pylori* which can cause adenocarcinoma of the stomach [13, 14]. The carcinogenic mechanism at play here is thought to be an entirely different one compared to carcinogenic viruses, namely sustained inflammation caused by a chronic, mostly subclinical infection [1]. For some links between infections and cancer, preliminary evidence has been presented, but the

Correspondence: sven.borchmann@uk-koeln.de

¹Department I of Internal Medicine, Center for Integrated Oncology Aachen Bonn Cologne Duesseldorf, University of Cologne, Cologne, Germany

²Cancer Center Cologne Essen – Partner Site Cologne, CIO Cologne, University of Cologne, Cologne, Germany

Full list of author information is available at the end of the article



© The Author(s). 2021 **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 licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence 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 licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

simultaneous presence of contradictory findings has led to widespread debate. An example of this is the finding that high levels of *Fusobacterium nucleatum* can be found throughout the cancerous tissue of colorectal cancer at much higher levels than in the tissue of benign adenomas or healthy colon mucosa [15–17]. Given the diversity of carcinogenic mechanisms [18], it is likely that other carcinogenic viruses and bacteria exist, although currently unknown.

Recent advances in massively parallel sequencing have made it possible to generate large amounts of data informing about the genome, transcriptome, and epigenome of a tissue [19]. Resulting datasets contain traces of non-host origin that are present either because of genomic integration or the presence of the virus or bacteria in the tissue itself. While some studies have already been performed with the goal of repurposing this data in order to reveal novel links between infections and cancer [4, 5, 20, 21], these resources have so far been underutilized.

With the above aim in mind, the present study leverages a large, high-quality collection of over 3000 whole genome sequencing datasets in order to gain insight into novel links between viruses, bacteria, and cancer.

Results

Samples

A total of 3025 whole genome sequencing datasets comprising 3.79 trillion reads were included in this study (Fig. 1a, Supplementary Table 1). These include 1330 whole genome sequencing datasets of tumor tissue samples across 14 different cancers from 19 International Cancer Genome Consortium (ICGC) [22] studies (Supplementary Table 1). Included patients were predominantly male ($n = 1028$, 60.6%) and elderly, with 47.3% of patients ($n = 801$) at least 60 years old (Fig. 1b,c). Two types of controls were used throughout this study. First, patient-matched normal (e.g., non-cancerous) tissue controls ($n = 1330$, mostly blood-derived or from tissue adjacent to the tumor, details in Supplementary Table 1) were utilized as controls. Only patients, for whom such a same-patient control was available were included in this study. Additionally, whole genome sequencing datasets of blood-derived DNA of 365 subjects from the 1000 genome project [23] were selected as a healthy control group substituting for the lack of negative sequencing controls to examine non-human DNA in the blood of healthy donors. Samples in the healthy control group were processed and sequenced at 5 different sequencing centers (86 at the BGI-Shenzhen, 86 at the Broad Institute, 11 at Illumina, 113 at the Sanger Institute, and 69 at Washington University in St Louis). Only subjects, in whom blood-derived DNA was directly subjected to whole genome sequencing, as opposed to DNA derived from immortalized

lymphoblastoid cell lines (LCL) (subset analyzed, $n = 102$), were included in the healthy control cohort. Whole genome sequencing datasets derived from LCL DNA showed a markedly different species-level taxon distribution likely representing taxa present due to LCL culture and not present in the donor itself (Supplementary Figure 1A, Supplementary Data 1-2). For validation purposes and to assess differential gene expression in cancers linked to certain species-level taxa, all available, matching RNA-seq datasets were also analyzed ($n = 324$).

Validation of pipeline

To perform taxonomic binning, a pipeline was built around Kraken [24], which at its core is based on the exact alignments of k-mers to their least common ancestor (LCA).

Kraken has been evaluated in two studies comparing metagenomic classifiers. In the first study [25], Kraken with its built-in filter performed well in species-level taxonomic binning, only being outperformed by few other tools, measured by F1 score, precision, recall, and area under the precision-recall curve. Importantly, it was only outperformed by tools that have a low recall if only very little sequence coverage is present for a taxon. These alternative tools are therefore not useful for this study. The same applies to combining different tools to arrive at a consensus binning. All tested combinations of other tools with Kraken have a very low recall rate at low coverage.

A low false positive rate is essential for this study and Kraken achieves this if its built-in filter is used. These performance characteristics were largely similar in another comparative study of metagenomic classifiers [26]; however, this, second study also found that the taxonomic binning by Kraken often contains false positive, very small bins. This can be mitigated by ignoring taxonomic bins that are very small, i.e., contain very few reads, which has been done in this study. Additionally, the filter stringency was increased. In one of the comparative studies of metagenomic classifiers, a filter threshold of 0.2 was used [25]. In the present study, a more stringent filter threshold of 0.5 was applied, aiming to further reduce false positives. Furthermore, the pipeline was validated twofold. In brief, it was confirmed that (i) the pipeline was able to identify already known bacterial and viral taxa in tissue-derived bacterial isolates and cell lines with known integration of viral DNA (Supplementary Data 2, Supplementary Figure 1B), and (ii) identified taxa in pairs of RNA-seq and WGS data of the same tumor tissue sample are correlated both in a combined dataset of all pairs and within each sample for which RNA-seq and WGS data was available ($n = 324$) (Supplementary Figure 1 C-D).

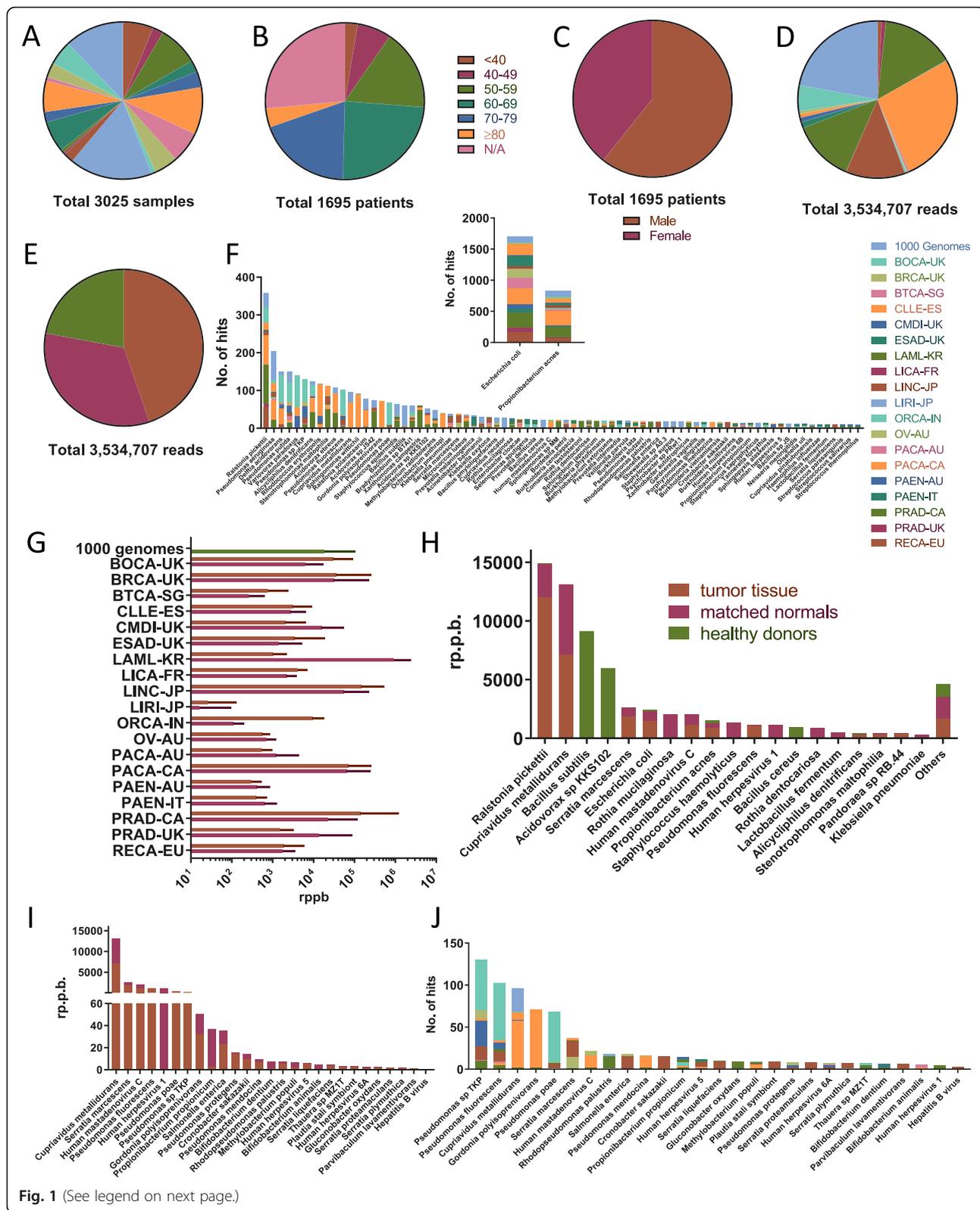


Fig. 1 (See legend on next page.)

(See figure on previous page.)

Fig. 1 Overview of sample characteristics and identified taxa. **a** Sample distribution by project. **b** Sample distribution by age group. **c** Sample distribution by gender. **d** Distribution of read pairs matching any species-level taxa by project. **e** Distribution of read pairs matching any species-level taxa by sample type. **f** Distribution of total reads per sample, analyzed unmapped reads per sample, and non-*phiX174* taxon-mapped reads per sample (non-*phiX174* taxon-mapped reads per sample not shown for 365 samples as none were detected). **g** Species-level taxa detected in at least 10 samples color-coded by project. **h** Detection of *Escherichia coli* and *Propionibacterium acnes* color-coded by project. **i** Average read pairs per billion (RPPB) detected across all species-level taxa by project and sample type. Bars show mean of samples and error bars show standard error of mean. **j** Average RPPB detected by species-level taxa and sample type. Average RPPB detected for all filtered species-level taxa identified as likely tumor-linked color-coded by project. *, tumor tissue samples were all primary solid tumor biopsy material for all projects with the following exceptions: CLL-ES, where tumor tissue samples were blood-derived CLL cells; CMDI-UK, where 3 tumor tissue samples were bone-marrow-derived and 29 samples blood-derived; LAML-KR, where all 12 tumor tissue samples were bone-marrow-derived; PACA-AU, where 1 tumor tissue sample was cell-line-derived; and 96 solid tumor biopsy material and PRAD-UK, where 4 samples were cancerous lymph nodes, 1 sample was a metastatic lesion, and 28 samples were solid tumor biopsy material.**, matched normal samples were all blood-derived for CLL-ES, LINC-JP, ORCA-IN, PRAD-CA, and RECA-EU, all matched non-cancerous tissue derived for BTCA-SG, LICA-FR, and PAEN-IT, a mix between blood-derived and matched non-cancerous tissue for BOCA-UK (69 vs. 7), BRCA-UK (44 vs. 1), ESAD-UK (87 vs. 10), LIRI-JP (250 vs. 6), PACA-AU (3 vs. 94), PACA-CA (55 vs. 68), and PAEN-AU (5 vs. 44), a mix between blood-derived and buccal cell-derived for CMDI-UK (31 vs. 1), a mix between blood-derived and EBV immortalized cell-line-derived for OV-AU (59 vs. 14), and a mix between blood-derived and cancer-free lymph node derived for PRAD-UK (23 vs. 10). For some samples, the source tissue was not specified (BRCA-UK ($n = 1$), LAML-KR ($n = 1$), and PACA-CA ($n = 24$)). ***, all healthy donor samples are blood-derived

A map of cancer-linked bacterial and viral taxa

A total of 3,534,707 read pairs matching bacterial, viral, or phage species-level taxa were detected across 19 studies in 3025 samples (Fig. 1d,e, Supplementary Table 1, Supplementary Data 4). Subsampling 10% of all read pairs did not alter the detected species-level taxa, and their relative composition compared to analyzing all non-human read pairs which was validated in a subset of patients (Supplementary Figure 1E-F, Supplementary Data 5). On average, 2.2 species-level taxa per sample were detected, although variation was high (Supplementary Data 6, Supplementary Figure 3 A-D). The mean number of total reads per sample was 1.25×10^9 (1.16×10^7 standard error of mean (s.e.m.)), the mean number of non-human, unmapped reads used for analysis per sample was 9.91×10^5 (4.34×10^4 s.e.m.), and the mean number of reads matching any taxon (excluding *phiX174*) was 2072 (291.8 s.e.m.) (Fig. 1f). A total of 218 species-level taxa could be identified in all examined samples (Fig. 1g, Supplementary Data 7). *Escherichia coli* and *Propionibacterium acnes* were the most detected species in all samples (Fig. 1h).

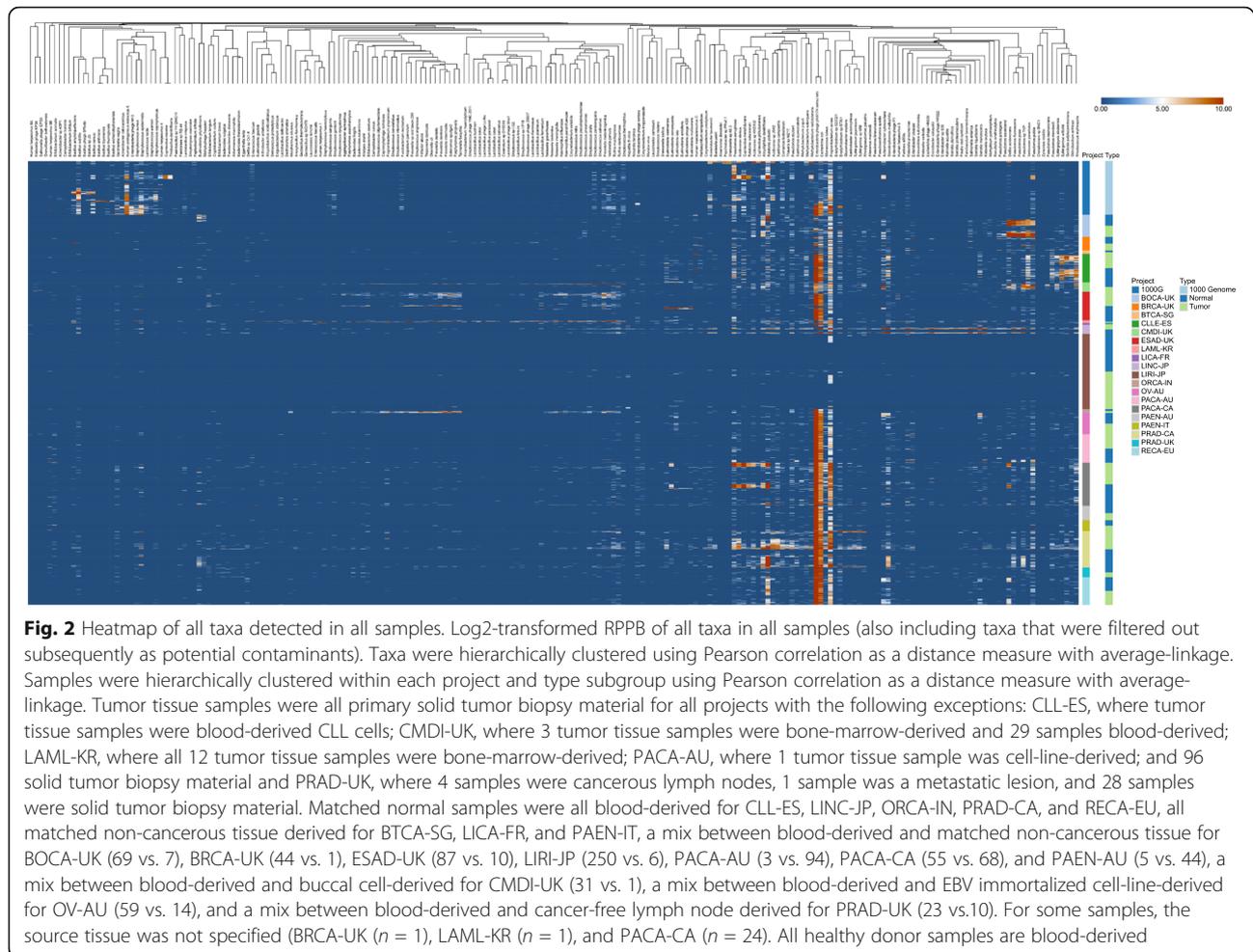
In order to control for differences in sequencing depth between samples, all raw read pair counts were normalized by dividing them by 1,000,000,000 total read pairs. This normalized count was defined as read pairs per billion (RPPB). The mean RPPB detected in healthy control samples from the 1000 genome cohort, matched normal samples, and tumor tissue samples were 18,112 (4238 s.e.m.), 20,003 (4295 s.e.m.), and 28,282 (9081 s.e.m.), respectively, with large variation across samples and sequencing projects (Fig. 1i). Of note, particularly high RPPB were detected in matched normal, saliva-derived samples from acute myeloid leukemia (AML) patients, as would be expected from a non-sterile source such as saliva. *Ralstonia pickettii* was the species with the

highest RPPB and almost absent in healthy donors, while being detected at higher levels in tumor tissue compared to matched normal samples (Fig. 1j). Except for *Bacillus subtilis*, *Acidovorax sp. KKS102*, and *Bacillus cereus*, all species with very high RPPB were detected predominantly in tumor tissue or matched normal samples. A clustered heatmap of the species-level taxa detected in all samples is provided in Fig. 2.

Filtering strategy

Next, filtering was performed to exclude taxa that were (i) frequently detected in the healthy control group, (ii) detected in only very few (< 5) tumor tissue or matched normal samples, (iii) phages, (iv) taxa that are commonly detected as part of the normal oral microbiome and were mainly detected in saliva, oral or esophageal cancer tissue samples (Supplementary Data 8), (v) taxa that have been previously described as sequencing contaminants (Supplementary Data 9), and (vi) taxa, for which the detected reads were unevenly distributed across the genome of the respective taxon (Supplementary Figure 4). After all these filtering steps (Supplementary Figure 5), 27 species-level taxa remained for further analysis (Supplementary Figure 6, Supplementary Data 10).

Among these, known tumor-linked taxa, such as *Hepatitis B virus* [10, 11] or *Salmonella enterica* [27, 28], were detected. Furthermore, taxa that have previously been implicated in carcinogenesis, although without enough evidence to support a carcinogenic role, such as *Pseudomonas* species [29, 30], and taxa that have never been implicated in carcinogenesis before, such as *Gordonia polyisoprenivorans*, were detected. Of note, most taxa in the final filtered species list were detected at much higher RPPB levels in tumor tissue compared to matched normal samples. However, lower RPPB of the respective species could still be detected in most



matched normal samples. Across all taxa, tumor tissue samples and matched normal samples were highly correlated (Supplementary Figure 1G).

Pan-cancer analysis

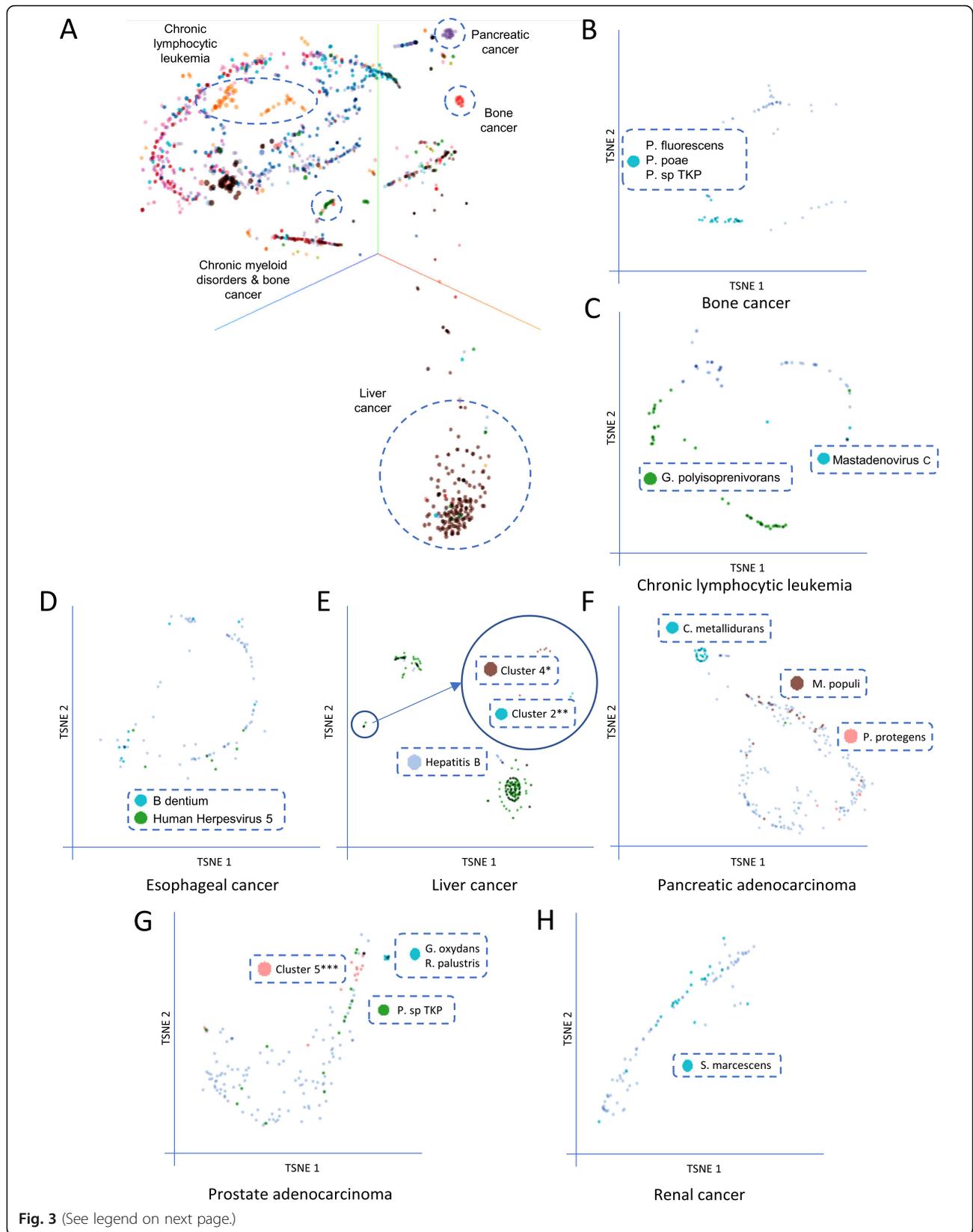
In order to better understand the link between detected taxa and different cancers and their relation to each other, a heuristic approach combining the non-linear dimensionality reduction and visualization method *t*-distributed stochastic neighborhood embedding (*t*-SNE) with *k*-means clustering was used. To focus on taxa that potentially play a more direct role in carcinogenesis and considering that RPPB of detected taxa were highly correlated between tumor tissue and matched normal tissue (Supplementary Figure 1G), only tumor tissues were included in the following analyses.

Utilizing a pan-cancer approach, a combined dataset of all tumor tissue samples was visualized using *t*-SNE using the log₂-transformed RPPB of all detected non-phage taxa ($n = 204$) as input variables. First, two distinct groups of patients with chronic lymphocytic

leukemia (CLL) could be identified. One group ($n = 24$) with detection of *Gordonia polyisoprenivorans* and another group ($n = 15$) with detection of *Pseudomonas mendocina* in the tumor tissue. Second, one distinct group ($n = 23$) of pancreatic cancer patients with detection of *Cupriavidus metallidurans* was observed. Third, one group ($n = 22$) of bone cancer patients with detection of *Pseudomonas poae*, *Pseudomonas fluorescens*, and *Pseudomonas sp. TKP* could be distinguished. Fourth, a group of patients ($n = 14$) with bone cancer ($n = 4$) or chronic myeloid disorders ($n = 10$) with detection of *Pseudomonas sp. TKP* was identified (Fig. 3a). Of note, all taxa that defined clusters in this analysis were detected across cohorts with the same and cohorts with different cancers making contamination unlikely.

Cancer-specific analysis

In order to gain further insight into links between detected taxa and specific cancers, each cancer type (Supplementary Data 11) was also analyzed separately. Patient groupings with similar detected taxa were visualized using *t*-SNE with log₂-transformed RPPB of all



(See figure on previous page.)

Fig. 3 Patient clusters defined by detected taxa can be identified across all patients and in cancer-type subgroups. **a** t-SNE visualization of all tumor tissue samples color coded by project using the \log_2 -transformed RPPB of all detected non-phage taxa as input variables. **b–h** t-SNE visualizations of single cancer subgroup tumor tissues (**b** bone cancer, **c** chronic lymphocytic leukemia, **d** esophageal cancer, **e** liver cancer, **f** pancreatic adenocarcinoma, **g** prostate adenocarcinoma, **h** renal cancer). Different colors indicate the *k*-means clusters of samples to indicate if *k*-means clustering and t-SNE visualization results in comparable sample groupings (Supplementary Figure 4). \log_2 -transformed RPPB of all detected non-phage taxa were used as input variables for t-SNE visualization. Dimension 1 and dimension 2 is shown on the x- and y-axis, respectively. Cluster 4*, *Pseudomonas sp.*, *Serratia sp.*, and *Salmonella enterica*; Cluster 2**, *Pseudomonas sp.*, *Serratia sp.*, *Salmonella enterica*, *Parvibaculum lavamentivorans*, and *Human Herpesvirus 5*; Cluster 5***, *Thauera sp.* MZ1T, *Cupriavidus metallidurans*, and *Pseudomonas mendocina*

species in the final filtered taxon list as input variables (Supplementary Data 10). Additionally, *k*-means clustering of patients was performed, using the same input variables. To confirm associations and clusters, two statistical tests were performed. First, sample RPPBs of a taxon in an identified cluster were compared to RPPBs of all samples of the same cancer type not in that cluster by Mann-Whitney *U* test to confirm an abundance difference of samples in the respective cluster compared to samples not in that cluster (denoted p_{cluster}). Second, sample RPPBs of a taxon in an identified cluster were compared to RPPBs of all samples in the healthy control cohort by Mann-Whitney *U* test to confirm the association of this taxon with a particular cancer (denoted p_{control}).

Furthermore, each cluster was linked with age, survival, gender, number of somatic mutations in known cancer genes, or specific somatic mutations in one of those cancer genes. Aiming to increase the validity of findings, all identified links between certain taxa and cancers were evaluated for presence in multiple independent sample cohorts of the same cancer, wherever possible. Multiple sample cohorts were available for liver cancer ($n = 3$), prostate adenocarcinoma ($n = 2$), pancreatic adenocarcinoma ($n = 2$), and pancreatic endocrine neoplasms ($n = 2$).

Bone cancer

In bone cancer, this dual methodology revealed a cluster of patients with detection of *Pseudomonas fluorescens* ($p_{\text{cluster}} = 3.2 \times 10^{-14}$, $p_{\text{control}} < 1 \times 10^{-15}$), *Pseudomonas sp. TKP* ($p_{\text{cluster}} = 1 < 10^{-15}$, $p_{\text{control}} < 1 \times 10^{-15}$), and *Pseudomonas poae* ($p_{\text{cluster}} = 2.8 \times 10^{-13}$, $p_{\text{control}} < 1 \times 10^{-15}$) (cluster 2), confirming the results of the pan-cancer analysis (Figs. 3 and 4b, a).

Chronic lymphocytic leukemia

In chronic lymphocytic leukemia, 2 taxon-linked clusters could be identified, one of patients with detection of *Human Mastadenovirus C* ($p_{\text{cluster}} < 1 \times 10^{-15}$, $p_{\text{control}} < 1 \times 10^{-15}$) (cluster 2) and one of patients with detection of *Gordonia polyisoprenivorans* ($p_{\text{cluster}} = 1.5 \times 10^{-10}$, $p_{\text{control}} = 5 \times 10^{-15}$) (cluster 3). Clusters could be identified by both methods, t-SNE and *k*-means (Fig. 3c,

Supplementary Figure 3B). Of note, the clusters of patients linked to *Human Mastadenovirus C* (cluster 2) and *Gordonia polyisoprenivorans* (cluster 3) were mutually exclusive ($p = 0.0124$) (Figs. 3 and 4c, b). There was a tendency towards different ages at diagnosis between the clusters ($p = 0.0743$) (Fig. 5a), with patients in cluster 2 (*Human Mastadenovirus C*) being younger. Additionally, there was a tendency towards a difference in survival between the different clusters ($p = 0.0745$). Patients not in any taxon-linked cluster had worse survival than patients in cluster 2 (*Human Mastadenovirus C*) or 3 (*Gordonia polyisoprenivorans*) ($p = 0.0246$) (Fig. 5b). Patients linked to *Gordonia polyisoprenivorans* (cluster 3) were more likely to have Binet C stage disease (5/36 vs. 1/61, Binet C vs. not, $p = 0.0252$) (Fig. 5c). These patients were also more likely to have TP53 mutations ($p(\text{cluster 3 vs. other}) = 0.0335$) (Fig. 5d).

Esophageal cancer

In esophageal cancer patients, 2 taxon-linked clusters were identified. *k*-means clustering revealed one cluster of patients with detection of *Bifidobacterium dentium* ($p_{\text{cluster}} = 1.7 \times 10^{-14}$, $p_{\text{control}} < 1 \times 10^{-15}$) (cluster 2) and one cluster with detection of *Human Herpesvirus 5* ($p_{\text{cluster}} = 7.8 \times 10^{-11}$, $p_{\text{control}} = 5 \times 10^{-15}$) (cluster 3). However, these 2 clusters could not be differentiated by t-SNE (Figs. 3 and 4d, c). There was a tendency towards different numbers of somatic mutations in cancer genes between clusters ($p = 0.0858$), with patients in cluster 3 (*Human Herpesvirus 5*) having fewer mutations than other patients ($p = 0.0392$) (Fig. 5e). Additionally, there was a tendency towards a difference in survival between the different clusters ($p = 0.1$). Patients in cluster 2 (*Bifidobacterium dentium*) or 3 (*Human Herpesvirus 5*) had worse survival than patients not in any taxon-linked cluster ($p = 0.040928$) (Fig. 5f).

Liver cancer

Patients with liver cancer could be grouped into 3 taxon-linked clusters by both *k*-means clustering and t-SNE. One cluster was defined by detection of *Hepatitis B virus* ($p_{\text{cluster}} < 1 \times 10^{-15}$, $p_{\text{control}} < 1 \times 10^{-15}$) (cluster 1). A second cluster was defined by detection of mainly *Pseudomonas* (p_{cluster} between 1.6×10^{-6} and 4.5×10^{-6} ,

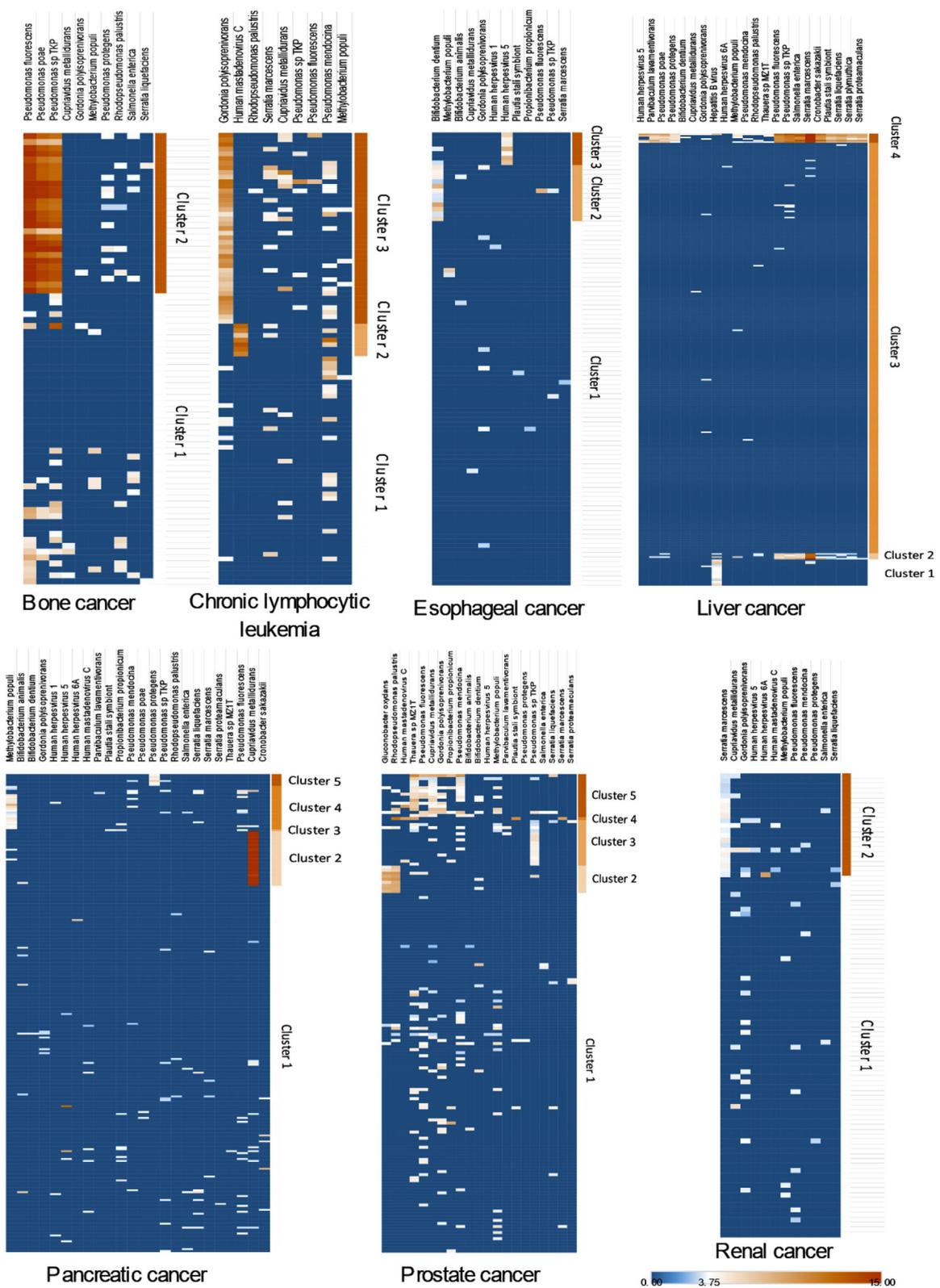


Fig. 4 Heatmaps of tumor-linked taxa for all cancers with discernible clusters. **a–g** \log_2 -transformed RPPB of all species-level taxa identified as likely tumor-linked and detected after filtering in all tumor-tissues of the indicated cancer. Results of *k*-means clustering of samples are shown

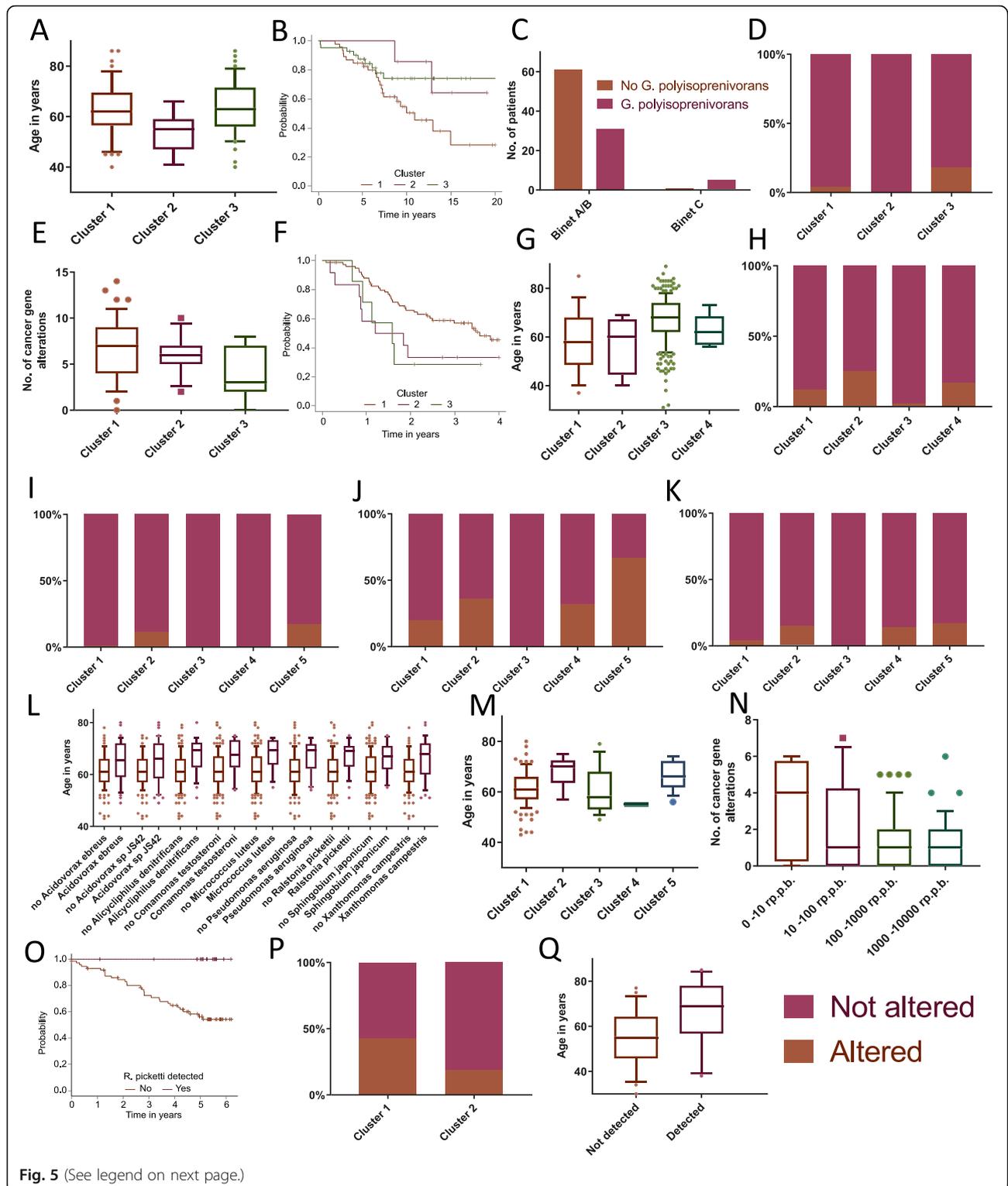


Fig. 5 (See legend on next page.)

(See figure on previous page.)

Fig. 5 Patient clusters defined by detected taxa are phenotypically distinct. **a** CLL patient clusters (1: no specific taxon link, 2: *Human Mastadenovirus C*, 3: *Gordonia polyisoprenivorans*) and age ($p = 0.0743$). **b** Survival by cluster in CLL ($p(1 \text{ vs. other}) = 0.0246$). **c** Binet stage by detection of *Gordonia polyisoprenivorans* ($p = 0.0252$). **d** TP53 mutation frequency by cluster in CLL ($p(\text{cluster 3 vs. other}) = 0.0335$). **e** Number of cancer consensus gene mutations by cluster (1: no specific taxon link, 2: *Bifidobacterium dentium*, 3: *Human Herpesvirus 5*) in esophageal cancer ($p = 0.0858$). **f** Kaplan-Meier survival curves for each cluster in esophageal cancer ($p(1 \text{ vs. other}) = 0.0409$). **g** Liver cancer patient clusters (1: *hepatitis B virus*, 2: *Pseudomonas sp.*, *Serratia sp.* and *Salmonella enterica*, 3: no specific taxon link, 4: *Parvibaculum lavamentivorans* and *Human Herpesvirus 5* in addition to taxa from cluster 2) and age ($p = 0.0015$). **h** RNF21 mutation frequency by cluster in liver cancer ($p = 0.0121$). **i** KMT2C mutation frequency by cluster (1: no specific taxon link, 2: *Cupriavidus metallidurans*, 3: no specific taxon link, 4: *Methylobacterium populi*, 5: *Pseudomonas protegens*) in pancreatic cancer ($p = 0.0308$). **j** CDKN2A mutation frequency by cluster in pancreatic cancer ($p = 0.0124$). **k** RNF21 mutation frequency by cluster in pancreatic cancer ($p(\text{RNF21}) = 0.0107$). **l** Detection of indicated taxa and age in prostate cancer (p_{adj} between 0.0015 and 0.0420). **m** Prostate cancer patient clusters (1: no specific taxon link, 2: *Gluconobacter oxydans* and *Rhodopseudomonas palustris*, 3: *Pseudomonas sp. TKP*, 4: no specific taxon link, 5: *Thauera sp. MZ1T*, *Cupriavidus metallidurans* and *Pseudomonas mendocina*) and age ($p = 0.0099$). **n** Link between *Propionibacterium acne* RPPB and number of cancer consensus gene mutations in prostate cancer ($p_{\text{adj}} = 0.0041$). **o** Kaplan-Meier survival analysis of *Ralstonia pickettii* detection status in renal cancer ($p_{\text{adj}} = 0.035$). **p** PBRM1 mutation frequency by cluster (1: no specific taxon link, 2: *Serratia marcescens*) in renal cancer ($p = 0.0723$). **q** *Pseudomonas sp. TKP* detection and age in chronic myeloid dysplasia ($p_{\text{adj}} = 0.039$). For all: The midline of the boxplots shows the median, the box borders show upper and lower quartiles, the whiskers show 5th and 95th percentiles and the dots outliers. Only tumor tissues were included in this figure

$p_{\text{control}} 1.3 \times 10^{-9}$ for all) and *Serratia species* (p_{cluster} between 2.8×10^{-6} and 1.1×10^{-4} , p_{control} between 1.3×10^{-9} and 4.8×10^{-7}) and *Salmonella enterica* ($p_{\text{cluster}} = 7.0 \times 10^{-7}$, $p_{\text{control}} = 1.3 \times 10^{-9}$) (cluster 2). A third cluster was defined by additional detection of *Parvibaculum lavamentivorans* ($p_{\text{cluster}} = 1.2 \times 10^{-12}$, $p_{\text{control}} = 2.9 \times 10^{-13}$) and *Human Herpesvirus 5* ($p_{\text{cluster}} = 5.0 \times 10^{-8}$, $p_{\text{control}} = 1.9 \times 10^{-8}$) as well as two additional *Pseudomonas species*, *Pseudomonas poae* ($p_{\text{cluster}} = 1.2 \times 10^{-12}$, $p_{\text{control}} = 2.9 \times 10^{-13}$), and *Pseudomonas protegens* ($p_{\text{cluster}} = 1.2 \times 10^{-12}$, $p_{\text{control}} = 2.9 \times 10^{-13}$) (cluster 4), in addition to those detected in the previous cluster (Figs. 3 and 4e, d). *Hepatitis B virus* and *Serratia species* were detected in 2 independent cohorts, while abovementioned *Pseudomonas species* were detected in all 3 independent cohorts, making contamination unlikely. There was a tendency towards different ages at diagnosis between the clusters ($p = 0.0015$) with patients in cluster 1 (*Hepatitis B*), cluster 2 (*Pseudomonas* and *Serratia sp.*), and cluster 4 (*Pseudomonas sp.*, *Serratia sp.*, *Parvibaculum lavamentivorans*, and *Human Herpesvirus 5*) being younger than patients not in any taxon-linked cluster ($p = 0.0001$) (Fig. 5g). In a similar pattern, these patients had a higher frequency of mutations in RNF21 ($p = 0.0121$) (Fig. 5h).

Pancreatic adenocarcinoma

In patients with pancreatic adenocarcinoma, 3 taxon-linked clusters could be identified using *k*-means clustering. One cluster was defined by detection of *Cupriavidus metallidurans* ($p_{\text{cluster}} < 1 \times 10^{-15}$, $p_{\text{control}} < 1 \times 10^{-15}$) (cluster 2), one by detection of *Methylobacterium populi* ($p_{\text{cluster}} < 1 \times 10^{-15}$, $p_{\text{control}} < 1 \times 10^{-15}$) (cluster 4), and a last one by detection of *Pseudomonas protegens* ($p_{\text{cluster}} = 3.6 \times 10^{-12}$, $p_{\text{control}} = 2.9 \times 10^{-13}$) (cluster 5). Only the cluster of patients linked to *Cupriavidus metallidurans* (cluster 2) could also be identified using t-SNE (Figs. 3

and 4f, e). All 3 taxon-species associations could be detected in both independent sample cohorts that were included, making contamination unlikely. Differences in the frequency of mutations between patients in different clusters were observed in KMT2C, CDKN2A, and RNF21. Patients in cluster 2 (*Cupriavidus metallidurans*), 4 (*Methylobacterium populi*), or 5 (*Pseudomonas protegens*) had a higher frequency of mutations in these genes compared to the majority of patients not in any taxon-linked cluster ($n = 187$) (p (KMT2C) = 0.0308, p (CDKN2A) = 0.0124, p (RNF21) = 0.0107) (Fig. 5i–k).

Prostate cancer

In patients with prostate cancer, 3 taxon-linked clusters emerged using *k*-means clustering, one defined by detection of *Gluconobacter oxydans* ($p_{\text{cluster}} = 6.6 \times 10^{-13}$, $p_{\text{control}} < 1 \times 10^{-15}$) and *Rhodopseudomonas palustris* ($p_{\text{cluster}} = 8 \times 10^{-15}$, $p_{\text{control}} < 1 \times 10^{-15}$) (cluster 2), one by detection of *Pseudomonas sp. TKP* ($p_{\text{cluster}} < 1 \times 10^{-15}$, $p_{\text{control}} < 1 \times 10^{-15}$) (cluster 3) and a last one by detection of *Thauera sp. MZ1T* ($p_{\text{cluster}} = 1.3 \times 10^{-11}$, $p_{\text{control}} < 1 \times 10^{-15}$), *Cupriavidus metallidurans* ($p_{\text{cluster}} < 1 \times 10^{-15}$, $p_{\text{control}} < 1 \times 10^{-15}$), and *Pseudomonas mendocina* ($p_{\text{cluster}} = 1.8 \times 10^{-8}$, $p_{\text{control}} < 1 \times 10^{-15}$) (cluster 5). Out of these, the cluster defined by *Gluconobacter oxydans* and *Rhodopseudomonas palustris* (cluster 2) and the one defined by *Thauera sp. MZ1T*, *Cupriavidus metallidurans*, and *Pseudomonas mendocina* (cluster 5) could be confirmed using t-SNE (Figs. 3 and 4g, f). *Cupriavidus metallidurans*, *Rhodopseudomonas palustris*, and *Pseudomonas species* could be detected in both independent sample cohorts that were included. However, *Gluconobacter oxydans* could only be detected in one cohort. There were age differences at diagnosis observable between the clusters ($p = 0.0099$) with patients in clusters 2 (*Gluconobacter oxydans*) and 5

(*Rhodopseudomonas palustris*) being older than the other patients ($p = 0.0005$) (Fig. 5m).

Renal cancer

In patients with renal cancer, a cluster of patients linked to *Serratia marcescens* ($p_{\text{cluster}} < 1 \times 10^{-15}$, $p_{\text{control}} < 1 \times 10^{-15}$) (cluster 2) was identified using *k*-means clustering, although t-SNE did not separate this group of patients (Figs. 3 and 4h, g). There was a tendency towards a lower frequency of PBRM1 mutations in patients in cluster 2 (*Serratia marcescens*) ($p = 0.0723$) (Fig. 5p).

Other cancers

In patients with pancreatic endocrine neoplasms, ovarian cancer, chronic myeloid disorders, and breast cancer, no discernable taxon-linked clusters could be identified (Supplementary Figure 7A-D).

Unbiased linkage analysis between bacterial and viral taxa and patient or cancer phenotypes

In addition to linking the above identified clusters with patient or cancer phenotypes, an unbiased analysis of links between the detection of a species-level taxa and patient or cancer phenotypes, such as age, survival, gender, number of somatic mutations in known cancer genes, and specific somatic mutations in one of those cancer genes, was performed utilizing both a pan-cancer approach and by analyzing each cancer type separately. In this analysis, all non-phage taxa detected ($n = 204$) were included and multiple testing correction was performed.

When analyzing cancer types separately, several links were identified. A group of bacterial taxa was linked to older patients in prostate cancer (p_{adj} between 0.0015 and 0.0420) (Fig. 5l). In chronic myeloid dysplasia, detection of *Pseudomonas sp. TKP* was also linked to older age ($p_{\text{adj}} = 0.039$) (Fig. 5q).

In the cancer-specific analysis, detection of *Ralstonia pickettii* was linked to improved survival in renal cancer, in fact no patients died ($p_{\text{adj}} = 0.035$) (Fig. 5o).

In prostate cancer, detection of and increasing *Propionibacterium acne* RPPB were linked to a decreasing number of cancer gene mutations ($p_{\text{adj}} = 0.0041$) (Fig. 5n).

Somatic lateral gene transfer

The integration of viral nucleic acids into the human host genome is well recognized as a carcinogenic process. High-level evidence exists for the integration of *Hepatitis B* [10, 11], *Human Papillomavirus* [8, 9], and *Epstein-Barr virus* (*Human Herpesvirus 4*) [31, 32], which are causally linked to hepatocellular carcinoma, cervical cancer, and lymphoma, respectively. Additionally, some evidence for somatic lateral gene transfer

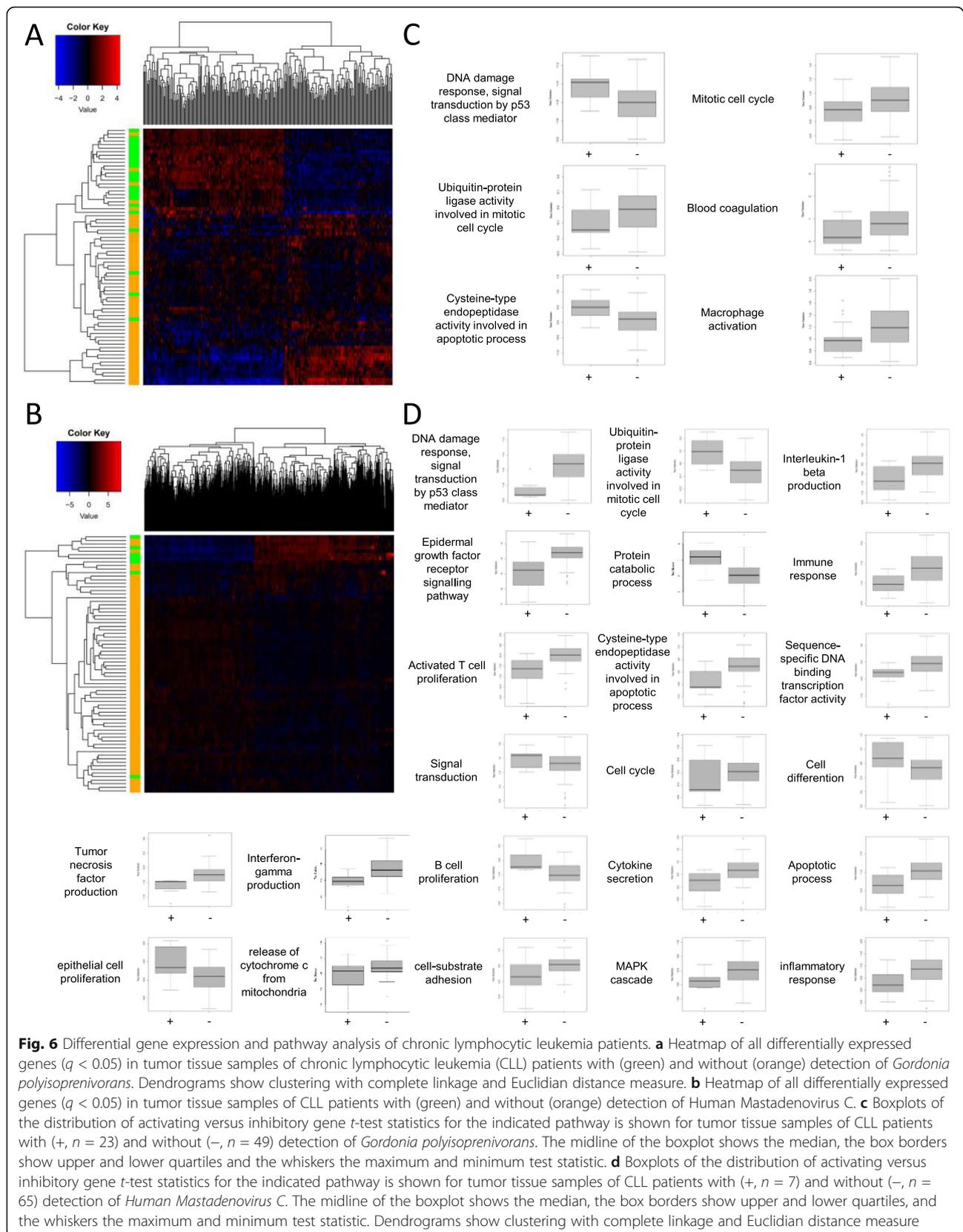
from bacteria to cancerous tissue has already been presented [20] and subsequently controversially discussed.

Aiming to find evidence for bacterial or viral DNA integration into the human host genome in this dataset, a pipeline for this purpose was developed. In brief, read pairs in which one read mapped to the human genome and one read to one of the taxa in the final filtered taxon list ($n = 27$) were counted and then compared to the number of read pairs in which both reads mapped to the respective taxa. The number of divergently mapping read pairs divided by the number of complete pairs mapping that respective taxa was used as a measure for genomic integration and lateral gene transfer. This was done for the tumor tissue datasets and the matched normal datasets including all patients ($n = 79$) for which the pipeline revealed at least 1000 RPPB matching one of the taxa in the final filtered taxon list (Supplementary Data 12). The highest rate of integration was observed for *Hepatitis B* in tumor tissue (9.62%) (Supplementary Figure 8A, Supplementary Data 13). Across all analyzed taxa, putative integrations were more common in matched normal samples than in tumor tissue samples ($p = 0.0009$, Wilcoxon paired signed rank test) (Supplementary Figure 8B, Supplementary Data 13-14) and for viral taxa compared to bacterial taxa ($p = 0.0001$, Mann-Whitney *U* test) (Supplementary Figure 8C, Supplementary Data 13-14). In conclusion, there was no evidence for a general phenomenon of lateral gene transfer for any species in the final filtered taxon list, with the exception of *Hepatitis B*, for which integration into the cancer genome has been widely described [10, 11].

Differential gene expression analysis

Differential gene expression analysis was performed for all taxa and cancer combinations with available tumor tissue RNA-seq data, in which at least 5 patients had upwards of 100 RPPB matching the respective taxon (Supplementary Data 15). Differentially expressed genes ($q < 0.05$) were identified for chronic lymphocytic leukemia patients with or without detection of *Gordonia polyisoprenivorans* (258 genes) (Fig. 6a, Supplementary Data 16), *Human Mastadenovirus C* (1725 genes) (Fig. 6b, Supplementary Data 17), and *Pseudomonas aeruginosa* (50 genes, Supplementary Data 18), respectively. Furthermore, differentially expressed genes were identified for ovarian cancer patients with or without the detection of *Escherichia coli* (22 genes) (Supplementary Data 19) and pancreatic adenocarcinoma patients with or without detection of *Propionibacterium acne* (3 genes) (Supplementary Data 20).

Next, differential gene expression data was used to perform bidirectional functional enrichment in order to identify pathways that are altered in patients with or without detection of the respective taxon. Only the two



patient groups with a significant number of differentially expressed genes, chronic lymphocytic leukemia with or without detection of *Gordonia polyisoprenivorans* (Supplementary Data 21) or *Human Mastadenovirus C* (Supplementary Data 22) had enriched pathways.

Overall, patients with detection of *Gordonia polyisoprenivorans* exhibited a gene expression pattern indicative of a decreased level of mitotic cell cycling, an increased DNA damage response, reduced blood coagulation, and reduced macrophage activation (Fig. 6c).

Contrary to that, patients with detection of *Human Mastadenovirus C* exhibited a gene expression pattern indicative of an increase in B cell proliferation, a reduction of tumor necrosis factor and interleukin beta 1 production, a reduction of activated T cell proliferation, and a decrease in cytokine secretion. In addition, the altered tumor tissue gene expression pattern of patients with detection of *Human Mastadenovirus C* indicated a markedly reduced DNA damage response (Fig. 6d).

Discussion

The aim of this study was to leverage a large, high-quality dataset of over 3000 samples to reveal novel links between viral and bacterial taxa and cancer. A total of 218 species-level taxa could be identified in tumor tissue, matched normal and healthy donor samples. Out of these, following extensive filtering, 27 taxa were likely cancer-linked. While studies of the viral metagenome of cancer tissues and patients have been performed using datasets from cancer genomics studies [4, 5, 21], similar large studies examining the bacterial metagenome are lacking.

Studies examining the viral metagenome of cancer tissues mainly identified known links of *Human Papillomavirus* to cervical and head and neck cancer, of *Hepatitis B* to liver cancer, and of *Human Herpesvirus 5* to a variety of cancers while the detection of *Human Mastadenovirus C* was controversial [4, 5, 21].

Smaller studies examining bacteria-tumor links in pan-cancer datasets have identified *Escherichia coli*, *Propionibacterium acne*, and *Ralstonia pickettii* in multiple cancers, while more specifically finding *Acinetobacter sp.* in AML and *Pseudomonas sp.* in both AML and adenocarcinoma of the stomach [20, 33]. Studying cancer-specific datasets, *Salmonella enterica*, *Ralstonia pickettii*, *Escherichia coli*, and *Pseudomonas sp.* were detected in breast cancer and adjacent tissue [34], while *Escherichia sp.*, *Propionibacterium sp.*, *Acinetobacter sp.*, and *Pseudomonas sp.* were frequently detected in prostate cancer [30]. Confirming the findings of the present study, these smaller studies found similar bacterial taxa, especially *Ralstonia pickettii*, *Escherichia coli*, *Propionibacterium acne*, *Salmonella enterica*, and *Pseudomonas sp.* Interestingly, the present study identified a number of taxa that

have not been previously identified in cancer tissue or matched normal samples of these patients, among them *Cupriavidus metallidurans*, *Gordonia polyisoprenivorans*, *Serratia sp.*, and *Bifidobacterium sp.*. Reasons for this are likely (1) the much higher number and diversity of patients and samples included, (2) the larger amounts of data examined per sample due to having higher-coverage WGS datasets available for all patients compared to RNA-seq or whole-exome sequencing (WXS) datasets in previous studies, and (3) the optimized and extensively validated bioinformatics approach used here. Of note, the filtering strategy used in this study to exclude taxa likely resulting from contamination has also eliminated bacterial taxa that have previously been linked to cancer such as *Escherichia coli* and *Propionibacterium acne*, mainly due to the frequent detection of these taxa in healthy donor samples from the 1000 genome cohort. Stringent filtering for species-level taxa that were detected with at least 100-fold higher RPPB in tumor tissue or matched normal tissue compared to healthy donor samples excluded these species. Despite that, *Propionibacterium acne* had about 1.4× and 3.1× higher RPPB for matched normal and tumor tissue, respectively, compared to healthy donors. These values were even higher in the case of *Escherichia coli*, namely 7.1× and 13.4× for matched normal and tumor tissue, respectively. Thus, it is well possible that tumor-linked taxa were eliminated due to the stringent filtering utilized. Therefore, some analyses in this study, such as the unbiased linkage analysis between detected taxa and phenotypes, were performed using the unfiltered dataset and all raw data is provided along with this manuscript for further analysis with different filtering approaches.

Examining specific cancer-pathogen links, a subgroup of bone cancer patients with detection of *Pseudomonas sp.* in the tumor tissue was identified. Consistent links between viral or bacterial taxa and bone cancer have not been described before, apart from some evidence linking *simian virus 40 (SV40)* infection to bone cancer [35]. Interestingly, *Pseudomonas sp.* are frequently implied in difficult to treat cases of osteomyelitis. Thus, one might speculate that chronic, subclinical infection exists and could be carcinogenic.

For CLL patients, two links were identified, one to *Gordonia polyisoprenivorans* and one to *Human Mastadenovirus C*. *Gordonia polyisoprenivorans* has not been linked to cancer before. Interestingly, the bacterium has been identified as a rare cause of bacteremia, so far exclusively in patients with hematological cancers [36–38]. Hematological cancers and their treatment are often associated with profound immunosuppression, allowing for infections with unusual environmental pathogens. It is conceivable that these patients had a latent infection with *Gordonia polyisoprenivorans* which exacerbated

into bacteremia and sepsis upon treatment-induced immunosuppression. To date, *Human Mastadenovirus C* has not been clearly linked to cancer, although recent reports have found it frequently detected in various cancer tissues, with one treating it as contamination [4]. Further indirect evidence for a role of both *Gordonia polyisoprenivorans* and *Human Mastadenovirus C* in CLL carcinogenesis is provided by (1) the observed age difference—patients with detection of *Human Mastadenovirus C* were markedly younger, (2) the mutual exclusivity of detection of *Gordonia polyisoprenivorans* and *Human Mastadenovirus C*, and (3) the fact that survival was different—patients with detection of either taxa had improved outcome. This was despite the higher likelihood of patients linked to *Gordonia polyisoprenivorans* having Binet C stage and patients linked to *Gordonia polyisoprenivorans* having a higher likelihood of having TP53 mutations, which are prognostically disadvantageous in CLL [39]. Strikingly, marked differences in host cancer tissue gene expression were observed for patients in which one of the taxa was detected with differences between cases with detection of *Gordonia polyisoprenivorans* and *Mastadenovirus C*. The observed tumor tissue gene expression pattern for patients linked to *Gordonia polyisoprenivorans*, especially an increased DNA damage response and reduced mitotic cell cycling, could explain the improved survival of these patients. The tumor tissue gene expression pattern of patients linked to *Human Mastadenovirus C* pointed towards reduced immune activity, especially reduced T cell function and a decrease in cytokine production and secretion, providing a potential explanation for the observed high detection frequency of *Human Mastadenovirus C* DNA in line with an uncontrolled infection due to a diminished immune response. Detection of *Mastadenovirus C* was linked to a decreased DNA damage response, which has been described as an important pathomechanism of CLL [40]. In addition to CLL patients, both *Gordonia polyisoprenivorans* and *Human Mastadenovirus C* were also detected frequently in other cancers, sometimes at high levels, while never being detected in healthy donors.

Two taxon-tumor links in esophageal cancer were identified, one to *Bifidobacterium dentium* and one to *Human Herpesvirus 5*. Interestingly, patients with detection of *Human Herpesvirus 5* had less somatic mutations in cancer consensus genes than the other patients, pointing to a possibly different carcinogenic mechanism, where the accumulation of multiple somatic mutations is not stringently needed for malignant transformation. Furthermore, patients not in any taxon-linked cluster had better survival. While the link of *Human Herpesvirus 5* to esophageal cancer is a novel finding, *Human Herpesvirus 5* has frequently been detected in adjacent

adenocarcinoma of the stomach [5]. Of note, overt *Human Herpesvirus 5* esophagitis can occur in immunocompromised hosts [41], which could be indicative of a latent infection of the esophagus by *Human Herpesvirus 5* in some hosts.

In patients with liver cancer, one subgroup of patients was defined by detection of *Hepatitis B*, a known cause of liver cancer [10, 11]. Interestingly, two other subgroups could be identified. *Pseudomonas sp.* and *Serratia sp.* were detected in both groups, with additional detection of *Parvibaculum lavamentivorans* and *Human Herpesvirus 5* in one group. There is some evidence that *Human Herpesvirus 5* might play a role in the carcinogenesis of liver cancer, among them detection of *Human Herpesvirus 5* DNA in tumor tissue and an increased seroprevalence in liver cancer patients [42] as well as frequent hepatitis in *Human Herpesvirus 5* infection underscoring hepatotropism of *Human Herpesvirus 5* [43]. The other identified taxa have not yet been implied in liver cancer carcinogenesis. Interestingly, hepatocellular carcinoma at younger age has been linked to chronic *Hepatitis B* infection [44]. Similarly, younger age of onset was also observed in liver cancers linked to the other taxa in this study.

Three bacteria-tumor links were identified in pancreatic adenocarcinoma, one with *Pseudomonas protegens*, one with *Methylobacterium populi* and one with *Cupriavidus metallidurans*. While *Pseudomonas sp.* have been shown to be a contributor to the pancreatic adenocarcinoma tissue microbiome [29], *Methylobacterium populi* and *Cupriavidus metallidurans* have not been detected in pancreatic cancer tissue. In fact, both taxa have not been discovered in human hosts but in environmental samples and are thus not considered part of the human microbiome, making it possible that they are contaminants not truly present in the tumor tissue samples analyzed. On the other hand, few infections of humans by these taxa have been described [45, 46] and, of note, the first published report of a *Cupriavidus metallidurans* infection was a case of septicemia in a patient with a pancreatic tumor [46].

In patients with prostate cancer, the most interesting findings were age differences between patient clusters defined by the detection of different taxa and a negative correlation of *Propionibacterium acne* detection and number of mutations in cancer consensus genes. Patients with detection of *Gluconobacter oxydans* and *Rhodopseudomonas palustris* as well as patients with detection of *Thauera sp. MZ1T*, *Cupriavidus metallidurans*, and *Pseudomonas mendocina* were markedly older than the other patients. *Gluconobacter sp.*, *Rhodopseudomonas sp.*, *Cupriavidus sp.*, and *Pseudomonas sp.* have previously been identified in prostate cancer and normal prostate tissue [30], but their precise role in prostate

disease is entirely unclear. *Propionibacterium acne* has been implied as a potential carcinogenic bacterium in prostate cancer [47, 48], possibly by creating a chronic inflammatory microenvironment [49]. In this study, the *Propionibacterium acne* detection frequency correlated negatively with the number of somatic mutations in cancer consensus genes. This could point to an alternative driver of carcinogenesis by chronic inflammation in the absence of accumulation of many mutations in cancer driver genes.

In renal cancer, a link to *Serratia marcescens* was identified in a subgroup of patients. While *Serratia marcescens* has been described as a frequent cause of urinary tract infections, especially in immunocompromised hosts in a nosocomial setting [50], it has so far not been implicated in cancer. Interestingly, survival of patients with detection of *Ralstonia pickettii* in their tumor tissue was markedly improved. *Ralstonia pickettii* is a bacterium that has been filtered out in this study because of frequent detection in healthy donor samples. It could be speculated that, while not causing any overt infection, low-level presence of *Ralstonia pickettii* in the human host is common and improves immunogenicity of renal cancer, thus, improving outcome. It has been shown that alterations of local and systemic immunity by the host microbiome influence the anticancer immune response [51], which might be highly relevant for a naturally immunogenic tumor, such as renal cancer [52].

The intriguing observation of increased somatic bacteria-human lateral gene transfer by Riley et al. [20] could not be made in this study. The only taxa for which more integration into the host genome was observed in tumor tissue compared to matched normal was *Hepatitis B*, for which integration into the genome of cancerous cells has been well recognized [53]. Additional integration of gut microbiome data would enhance this manuscript. The gut microbiome has recently emerged as being highly relevant to carcinogenesis, especially of cancers exposed to it, such as colorectal cancer [54]. It has also emerged that the gut microbiome can modulate cancer treatment efficacy, particularly of immunotherapy [55]. However, gut microbiome data was not available for the patients included in this study.

It is important to note, that this study is an explorative analysis of potential novel relationships between cancers, viruses, and bacteria and only experimental validation can really prove the postulated links of bacterial and viral taxa to certain cancers. Nevertheless, every attempt has been made to reduce false positives, by carefully choosing a pipeline based on available comparison studies [25, 26], applying stringent filtering, removing low complexity sequence, and removing taxonomic bins with only very few hits. Looking at genetic variation on the strain level between different patients could further

validate the findings and rule out contamination [56]. However, this is not possible, as methods to do so all require a minimum coverage of the taxon in question. Due to the nature of this study—examining low-level presence of viruses or bacteria in cancer tissues—sufficient coverage is not available. Further validation could come from validating some of the findings on long-read sequencing platforms [57]. However, this study reused datasets already available and had no access to the original samples; therefore, such a further validation is not possible.

Conclusions

In conclusion, the present study provides an unprecedented atlas of links of both bacterial and viral taxa to cancer. In addition to confirming known or recently postulated links, several novel associations between, bacteria, viruses, and cancer were identified across multiple cancer entities, laying the groundwork for further studies and experimental validation.

Methods

Data sources and data availability

Mapped sequencing data for the included ICGC studies was obtained via the ICGC DCC [22, 58] and downloaded using customized scripts. The use of controlled access ICGC data for this project was approved by the ICGC data access compliance office. Mapped sequencing data for the 1000 genome healthy control samples was obtained from the 1000 genome FTP server [59] using customized scripts. Sequencing data used for validation was obtained from the European Nucleotide Archive (ENA) [60]. Each sample is clearly identified by the identifiers in Supplementary Data 4. These identifiers can be used to obtain further sample and donor information from the ICGC data portal (<https://dcc.icgc.org/>). While basic data is available without approval, downloads of raw sequencing data via the ICGC data portal have to be requested from the ICGC data access compliance office. Raw read counts and analysis data are available in full and included in Supplementary Data 4.

Computing environment

Data analysis was performed using a HP Z4 workstation in a Unix environment either using software as mentioned throughout the “Methods” section or customized scripts. Some analyses were performed employing the Galaxy platform [61].

Pipeline for taxonomic classification

First, unmapped (non-human) read pairs were extracted from a random 10% subsample of each sample’s downloaded sequencing data using Samtools (version 1.7) [62]. Subsampling 10% did not alter the detected

species-level taxa and their relative composition (Supplementary Figure 1E-F, Supplementary Data 5). Bam files were sorted by query name with Picard Tools (version 2.7.1.1) [63] and converted to FastQ files using Bedtools (version 2.26.0.0) [64]. Subsequently, Trimmomatic (version 0.36.3) [65] was used to trim reads with sliding window trimming with an average base quality of 20 across 4 bases as cutoff and dropping resulting reads with a residual length < 50. Read pairs, in which one read was dropped according to these rules, were dropped altogether. Remaining read pairs were joined with FASTQ joiner (version 2.0.1) [66] and converted to FASTA files using the built-in function FASTQ to FASTA (version 1.0.0) from Galaxy [61]. Next, VSearch (version 1.9.7.0) [67] was used to mask repetitive sequences by replacing them with Ns using standard settings. These masked and joined read pairs were fed into Kraken (version 1.1.1) [24] using a database of all bacterial and viral genomes in Refseq (release 85). The output of each run was filtered with Kraken (version 1.1.1) [24] setting a confidence threshold of 0.5. A report combining the output of all samples and runs was generated using Kraken (version 1.1.1) [24]. The output was then arranged using customized scripts in R (beginning with version R 3.3.2. and subsequently updated) [68] to generate the raw metagenome output of each sample. Next, the raw metagenome output of each sample was filtered by only including species-level taxa and excluding all species-level taxa that were supported by less than 10 read pairs across all samples using R (beginning with version R 3.3.2. and subsequently updated) [68]. Read pairs assigned to Enterobacteria phage phiX174, which is ubiquitously used as a spike-in control in next-generation sequencing were omitted from all counts and analyses as an intended contaminant, except for Fig. 2, which aims to visualize the full, raw dataset.

To correct for the variation of sequencing depth across samples, matched read pairs per billion read pairs raw sequence (RPPB) were calculated for each sample and each taxon.

RPPB was calculated using the following formula:

$$10^9 \times \frac{\text{read pairs of a sample assigned to a given taxon}}{\text{total available read pairs for a given sample}}$$

A taxon was heuristically considered detected, when a respective sample had at least 100 RPPB assigned to that taxon.

Filtering strategy

First, all taxa that were also highly prevalent in the healthy control group were excluded from further analysis. In detail, a taxon was required to have a mean RPPB across either the tumor tissue samples or the

matched normal samples compared to the healthy control samples of at least 100-fold higher, to be included. This cutoff excluded all taxa that were also highly prevalent in the healthy control samples while at the same time allowing to further analyze taxa that were enriched in matched normal samples such as blood as well as taxa that were dominant in tumor tissue. After this step, 147/218 (67.4%) potential tumor-linked species-level taxa remained. Next, taxa that were detected in fewer than 5 tumor tissue or matched normal samples were excluded ($n = 78$), as well as all remaining phages ($n = 2$). *Hepatitis B* as a known cancer-linked virus was re-included despite being detected in fewer than 5 tumor tissue or matched normal samples according to these criteria so that 68/218 (31.2%) taxa remained.

Taxa that survived this filtering strategy were likely to be tumor-linked but could also represent artifacts from contamination. To account for that, all taxa were filtered out that are known to be regularly present in the oral microbiome [69–71] and were at the same time mainly (> 50% of all RPPB matching a respective taxon) detected in samples that are likely contaminated with the oral microbiome (saliva matched normal, oral cancer tissue, and esophageal cancer tissue samples) (Supplementary data 8). For example, this filtered out taxa that are a commonly present in the oral human flora such as *Streptococcus mitis* and were indeed mainly detected in tumor tissue samples of oral cancer or esophageal cancer and likely a contaminant based on the biopsy location. Similarly, these taxa were detected in the matched normal of leukemia cases whose matched normal was a saliva sample likely containing taxa of the normal human microbiome. After this filtering step, 49/218 (22.5%) species-level taxa remained.

It has recently emerged that both reagents and kits used in DNA extraction and library preparation as well as ultrapure water used in laboratories can contain contaminants that can hamper the detection of truly present taxa in low biomass or high background (e.g., human) samples. As this study was performed using samples that were both low in non-human biomass and in the context of high human background, the aim was to further reduce false positives by compiling a list of common contaminants in microbiome studies. Recommended approaches, such as the sequencing of blank controls [72], were not feasible as the present study was conducted utilizing already sequenced primary material. Therefore, a database of common contaminants from various studies [73–77] examining this issue was compiled (Supplementary data 9). All previously recognized contaminant taxa apart from those that were described as a contaminant on the genus level but where different species within the genus were detected differently in 1000 genome control samples and matched normal or tumor tissue

samples were excluded. This was the case for the genera *Pseudomonas* and *Methylobacterium*. While the species *Pseudomonas aeruginosa*, *Pseudomonas putida*, and *Pseudomonas stutzeri* were detected in both 1000 genome control and matched normal or tumor tissue samples and thus likely contaminants, *Pseudomonas fluorescens*, *Pseudomonas mendocina*, *Pseudomonas poae*, *Pseudomonas protegens*, and *Pseudomonas sp. TKP* were not detected in 1000 genome healthy control samples. Thus, it is likely that the species-level resolution of this analysis was able to differentiate between common contaminants and possibly tumor-linked taxa. Similarly, differentiation was possible between *Methylobacterium radiotolerans* and *Methylobacterium extorquens* (likely contaminants) and *Methylobacterium populi*, which was only found in tumor samples. The same held true for *Cupriavidus metallidurans* and *Propionibacterium propionicum*. Of note, the 1000 genome healthy control cohort used within this study contains samples processed and sequenced in 5 different sequencing centers (86 at the BGI-Shenzhen, 86 at the Broad Institute, 11 at Illumina, 113 at the Sanger Institute, and 69 at Washington University in St Louis). Thus, this control cohort itself serves as a bona fide contamination control including potential sequencing contaminants originating in different reagents, different suppliers, and different laboratory or environmental contaminants. After this filtering step, 27/218 (12.4%) species-level taxa remained (Fig. 1i,j, Supplementary Figure 6). While it cannot be ruled out that all these filtering steps removed truly cancer-linked taxa, the aim was to be cautious and rather accept a false-negative than a false positive finding. Of course, experimental validation of the relevance of one of the taxa eliminated by one of the filters could prove that this taxon is both, a sequencing contaminant and a relevant taxon in cancer.

If a taxon is indeed present in a tissue, matching read pairs are expected to be uniformly distributed across its genome. Consequently, a further filtering step was introduced. The sequencing data from all samples was combined and matched against a reference database constructed out of the genome of these 27 taxa. Next, the coverage distribution of reads across each taxon's genome was assessed (Supplementary Figure 4). If detection of the respective taxa is the result from misalignment or sequence similarity between the taxa and for example cloning vectors used in the production of sequencing reagents, an uneven coverage would result. It was found that all read pairs matching *Human Mastadenovirus C* aligned to short parts of its genome with a maximum length of a few hundred base pairs and abrupt drops in coverage (Supplementary Figure 4). This was also found in another study using different cancer tissue sequencing data and resulted in excluding

Mastadenovirus C from further analysis [4]. Using Blast (beginning with version 2.7.1 and subsequently updated) [78], it was found that read pairs matching *Human Mastadenovirus C* aligned all equally well to commonly used cloning vectors, such as pAxCALGL, which might have been used in the production of reagents used for sequencing. This form of contamination was recently analyzed and found to be frequent [79]. However, read pairs aligning to *Human Mastadenovirus C* originated from very few, seemingly unlinked samples from diverse cancer sequencing projects, making contamination by recombinant DNA unlikely. Another explanation for the observed coverage pattern is somatic genomic integration of parts of the *Human Mastadenovirus C* genome into a specific cancer genome. On balance, *Human Mastadenovirus C* was therefore not excluded from further analysis.

Clustering of pipeline hits

t-SNE was performed using a web-based TensorFlow Embedding Projector implementation [80]. The learning rate and the perplexity were heuristically set to 10 and 30 for all analyses, respectively, except for the liver cancer subset, for which the perplexity was set to 50 due to improved cluster discrimination. The number of iterations was heuristically chosen, so that no major changes of cluster composition occurred upon increasing the number of iterations. Depending on the subset analysis, between 500 and 1500 iterations were needed to reach that point. t-SNE was performed in 3 dimensions for the pan-cancer analysis and in 2 dimensions for the cancer-specific analyses.

k-means clustering was performed using Morpheus [81]. Unsupervised *k*-means clustering using Euclidean distance as a similarity measure was employed with the number of clusters being heuristically informed by combining visual inspection, comparing t-SNE and *k*-means clusters and by examining marginal reduction of within-group variance with increasing numbers of clusters (i.e., the elbow method) (Supplementary Figure 9 A-J).

To combine the information obtained by both complementary methods, t-SNE clustering was repeated with the same settings for each analysis, while color-coding clusters inferred from *k*-means clustering.

Assessment of taxonomic differences between cell culture-derived and blood-derived DNA samples from the 1000 genome project

All taxa in the final taxon list ($n = 218$) (Supplementary data 10) which were identified in blood-derived 1000 genome healthy control samples were selected ($n = 76$) (Supplementary data 1). The pipeline was subsequently applied to randomly selected (identifier ending with 8 or 8) LCL-derived 1000 genome samples ($n = 102$) (2).

Finally, RPPB in these LCL-derived samples were calculated for all 76 taxa that were identified in the blood-derived 1000 genome samples and compared between blood-derived and LCL-derived samples.

Assessment of effect of subsampling of read pairs on relative taxon distribution

To show that subsampling alters neither the detected species-level taxa nor their relative composition compared to analyzing all non-human read pairs, a subset of 184 tumor tissue samples (Supplementary data 6) was analyzed, without any subsampling and subjected to the pipeline in the same way as in the main analysis. For example, read pairs matching *Human Mastadenovirus C*, *Pseudomonas poae*, *Ralstonia pickettii*, and *Propionibacterium acnes* were used to compare absolute read pair counts between full data and the subsample for all taxon-sample pairs in which the 10% subsample had at least 10 matches to the respective taxon.

Concordance between WGS and RNA-seq

In order to assess the concordance between RNA-seq and WGS experiments performed on the same sample, the main pipeline was applied with the same settings to all samples with RNA-seq and WGS paired data available ($n = 324$). Pearson correlation coefficients of \log_{10} transformed data were calculated for both a combined dataset of all RNA-seq / WGS pairs and for each sample for which RNA-seq and WGS data was available ($n = 324$).

Assessment of somatic lateral gene transfer

In order to assess the integration of bacterial DNA into human DNA, read pairs with one read matching the human genome and one read matching one of the taxa in the final filtered taxon list ($n = 27$) (Supplementary data 10) were identified. First, representative bacterial and viral genomes for the final filtered taxon list were downloaded (Accession numbers in Supplementary data 10). These genomes were merged with the human reference genome (hg1k_v37, downloaded from the 1000 genome FTP server [59]) into one FASTA file using customized scripts. Second, all read pairs in which only one read of a read pair was mapped to the human genome were filtered using Samtools (version 1.7) [62]. All tumor tissue samples and matched normal samples of all patients in which at least 1000 RPPB matching one of the taxa in the final filtered taxon list were identified in that respective patient's tumor tissue sample ($n = 79$, Supplementary data 12) and were included in this analysis. BWA mem (version 0.7.17) [82] with standard settings in paired end mode was used to align all such read pairs to the merged FASTA file of all taxa and the human reference genome. Next, all read pairs that were now

divergently mapped were extracted using Samtools (version 1.7) [62] and customized scripts by only including read pairs where one read mapped to one of the included non-human taxa and the other read mapped to a human sequence. Only read pairs with a mapping quality of at least 40 were retained. Customized scripts were used to count and tabulate all obtained divergently mapped read pairs by taxa and sample, respectively (Supplementary data 13-14). In order to normalize read pairs mapping to putative integration sites (i.e., divergently mapped read pairs as defined above) by correcting for the total number of read pairs matching a taxon with a similar approach (i.e., non-divergently mapped, putatively non-integrated reads), a comparable pipeline was applied to the data used for the main analysis (i.e., both reads in a read pair not mapped to the human genome) of all patients included in the integration analysis ($n = 79$) (Supplementary data 12). First, the genomes of all taxa in the final filtered taxon list were downloaded (Accession numbers in Supplementary data 10) and merged without adding any further human sequences into one FASTA file to create a reference genome containing all taxa in the final filtered taxon list. BWA mem (version 0.7.17) [82] with standard settings in paired end mode was used to align these reads to the merged FASTA file of all taxa in the final filtered taxon list. Subsequently, Samtools (version 1.7) [62] was used to filter the aligned data to only include read pairs that mapped as a proper pair to only one taxon with a minimum mapping quality of 60. Samtools (version 1.7) [62] and customized scripts were used to count and tabulate all mapped reads. The integration rate of a taxon in either tumor tissue samples or matched normal samples was calculated by dividing the number of divergently mapping read pairs by the number of read pairs mapping as a proper pair to the respective taxon.

Assessment of links of taxon-defined patient clusters to patient or cancer phenotypes

Differences in age at diagnosis between patient clusters were first analyzed by ANOVA for each cancer. All results with $p_{\text{anova}} \leq 0.1$ are shown in Fig. 5. Such clusters or combinations of clusters were then compared to the other clusters by Student's t test.

Differences in the gender distribution between patient clusters were analyzed by chi-square test for each cancer.

In order to analyze relationships between the number or type of cancer-associated somatic mutations and detection of specific taxa, all somatic mutations for all patients included in this study were obtained from the ICGC data portal [22, 58]. All synonymous mutations were filtered out. Subsequently, only Tier 1 cancer gene census [83] genes that were altered in more than 20

cases were filtered using customized scripts and included in the analysis (Supplementary data 23).

Differences in the number of somatic mutations in cancer genes between patient clusters were first analyzed by ANOVA for each cancer. All results with $p_{\text{anova}} \leq 0.1$ are shown in Fig. 5. Such clusters or combinations of clusters were then compared to the other clusters by Student's t test.

The Kaplan-Meier method was used to estimate survival curves for each patient cluster in each cancer type with available survival data. Differences in survival between patient clusters were analyzed by log rank test.

Links between patient clusters and somatic mutations in single cancer genes were analyzed if the gene was altered in at least 10 patients in that respective cancer type. Clusters or combinations of clusters were then compared to the other clusters by Fisher's exact test.

All calculations were performed with R (beginning with version R 3.3.2. and subsequently updated) [68].

Unbiased linkage analysis between single bacterial and viral taxa and patient or cancer phenotypes

For this analysis, a taxon was considered detected in a patient if at least 100 RPPB matched the respective taxa in the patient's tumor tissue sample. All phages were excluded from the analysis. All included projects were grouped by cancer (Supplementary data 11). Cancer genes were defined as above.

All calculations were performed with R (beginning with version R 3.3.2. and subsequently updated) [68]. All p values were corrected for multiple testing using the FDR method to obtain a q -value, which was considered significant if < 0.05 .

Survival analysis

First, the Kaplan-Meier method was used to estimate survival curves for each cancer-taxon pair that was detected in at least 10 patients. Differences in survival between patients with or without detection of a respective taxon were analyzed using log rank tests, stratified by ICGC project. Additionally, a pan-cancer analysis was performed in the same way, also stratifying by ICGC project.

Links between bacterial or viral taxa and patient gender

Links between detection of a taxa and patient gender were analyzed by Fisher's exact test for each cancer-taxon pair that was detected in at least 10 patients. Additionally, a pan-cancer analysis was performed in the same way, using a logistic regression model that included the ICGC project as an independent variable.

Links between bacterial or viral taxa and patient age

Links between the detection of a taxon and patient age at diagnosis were analyzed by Student's t test for each cancer-taxon pair that was detected in at least 10 patients. Additionally, a pan-cancer analysis was performed in the same way, using a linear regression model that was stratified by ICGC project.

Links between bacterial or viral taxa and number of somatic mutations in cancer genes

Links between the detection of a taxon and the number of non-synonymous somatic mutations in cancer consensus genes [83] of a patient were analyzed by Student's t test for each cancer-taxon pair in which a taxon was detected in at least 10 patients. Additionally, a pan-cancer analysis was performed in the same way, using a linear regression model that was stratified by ICGC project.

Links between bacterial or viral taxa and specific somatic mutations

Links between the detection of a taxon and non-synonymous somatic mutations in one of the cancer consensus genes [83] of a patient were analyzed by Fisher exact test for each cancer-taxon pair that was detected in at least 10 patients. Additionally, a pan-cancer analysis was performed in the same way, using a logistic regression model that included the ICGC project as an independent variable.

Assessment of differential gene expression

Reads per kilobase of transcript, per million mapped reads (RPKM) data was downloaded from <http://dcc.icgc.org/pcawg> for all available tumor tissue samples. Ensemble gene ID was substituted by the standard Human Genome Organization (HUGO) Gene Nomenclature Committee (HGNC) symbol downloaded from <http://genenames.org/download/custom> and the RPKM data was then linked to the identified species-level taxa in each sample using R (beginning with version R 3.3.2. and subsequently updated) [68]. sRAP [84] was used to normalize RPKM values, perform quality control and differential gene expression analysis, and to identify pathways that are differentially expressed. For these analyses, each cancer was analyzed separately and patients that had more than 100 RPPB matching one species-level taxon in the final filtered taxon list were compared with those who did not. Differential gene expression analysis was performed for all taxa that were detected with more than 100 RPPB in at least five patients and RNA-seq data available (Supplementary data 15). Next, sRAP [84] was used to perform bidirectional functional enrichment of gene expression data to identify pathways up- or downregulated between patients with or without

detection of a taxon. Briefly, the distribution of activation versus inhibition *t*-test statistics for all samples linked or not linked to a specific taxon was compared using ANOVA and corrected for multiple testing using the FDR method [84]. A gene set was considered functionally enriched if $q < 0.05$. Gene sets likely not relevant for the respective cancer were excluded and gene sets provided with sRAP [84] were reduced to gene ontology (GO) gene sets [85].

Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s40168-021-01039-4>.

Additional file 1: Supplementary Table 1. Included patients and samples. **Supplementary Figure 2.** Validation of pipeline and analytical approach. A, mean RPPB detected for indicated species in blood-derived and lymphoblastoid cell line 1000 Genome samples sorted by mean of blood-derived samples. B, proportion of non-human read pairs matching the indicated taxon of read pairs matching any species-level taxon for each external validation sample. C, comparison of Kraken matched read pairs in RNA-seq and WGS data of the same sample. Each dot represents one species-level taxon in one sample with both RNA-seq and WGS data available. The line represents the best-fitted line (log-log linear regression). Pearson correlation coefficients (log-log) are shown with two-sided *p*-values. D, plot of Pearson correlation coefficient (log-log) distribution of all samples with both RNA-seq and WGS data available. Each dot represents the Pearson correlation coefficient within a single sample. E, comparison of 10% subsample and full dataset. Each dot represents one species-level taxon in one sample for which both the full and the subsampled dataset has been analyzed and indicates the absolute read count identified in both samples. The line represents the best-fitted line (log-log linear regression). Pearson correlation coefficients (log-log) are shown with two-sided *p*-values. F, Ratio of absolute read counts in the full sample to the 10% subsample for 4 selected taxa. The mean ratio of all samples in which the respective taxon was detected is indicated by the symbol and the error bars indicate the standard error of the mean. The dotted line shows the expected ratio of 10. G, comparison of tumor tissue and matched normal by patient and taxon. Each dot represents one species-level taxon in one patient with both tumor tissue and matched normal analyzed and indicates the RPPB in both samples. The line represents the best-fitted line (log-log linear regression). Pearson correlation coefficients (log-log) are shown with two-sided *p*-values. **Supplementary Figure 3.** Alpha diversity. A, counts of 1000 Genome samples by species-level richness. B, counts of tumor tissue samples by species-level richness color-coded by project. C, counts of matched normal samples by species-level richness color-coded by project. D, comparison of richness between projects and sample type. Bars show mean and error bars standard deviation. N in brackets indicates total sample number for each project. **Supplementary Figure 4.** Coverage distribution for all tumor-linked species-level taxa. Coverage distribution across each species-level taxon identified as tumor-linked. **Supplementary Figure 5.** Flow chart of taxa filtering strategy. Flow chart of filtering strategy to derive likely tumor-linked species-level taxa. **Supplementary Figure 6.** Heatmap of filtered taxa. Log₂-transformed RPPB of all species-level taxa identified as likely tumor-linked after filtering in all samples. Taxa were hierarchically clustered using Pearson correlation as a distance measure with average-linkage. Samples were hierarchically clustered within each project and type subgroup using Pearson correlation as a distance measure with average-linkage. **Supplementary Figure 7.** Heatmaps of tumor-linked taxa for all cancers without discernible clusters. A-D, log₂-transformed RPPB of all species-level taxa identified as tumor-linked and detected after filtering in all tumor-tissues of the indicated cancer. Results of *k*-means clustering of samples are shown. **Supplementary Figure 8.** Host integration. A, integration rate by species for tumor tissue and matched normal sample. B, difference in integration rates between bacterial and viral taxa ($p < 0.0001$, Wilcoxon rank-sum test, two-

tailed). The midline of the boxplot shows the median, the box borders show upper and lower quartile, the whiskers show 5th and 95th percentiles and the dots outliers of species-specific integration rates in tumor tissue or matched normal samples. C, difference in integration rate between tumor tissue and matched normal samples ($p = 0.0009$, Wilcoxon signed rank test, two-tailed). The midline of the boxplot shows the median, the box borders show upper and lower quartile, the whiskers show 5th and 95th percentiles and the dots outliers of species-specific integration rates. **Supplementary Figure 9.** "Elbow method" to determine *k* for *k*-means clustering. A-J, plot of reducing within group sum of squares for increasing *k* (number of clusters) in *k*-means clustering (Supplementary Figures 4 and 5) of log₂-transformed RPPB of all species-level taxa identified as tumor-linked and detected after filtering in all tumor-tissues for each indicated cancer.

Additional file 2: Supplementary 1-3. control and validation.

Supplement 4. rppb and read counts raw data_revised.

Supplement 5. downsampling vs full dataset. **Supplement 6.** number of species detected in each sample. **Supplement 7-9.** filtering and grouping. **Supplement 10.** final taxa hit list. **Supplement 11.** cancer study groupings. **Supplement 12-14.** integration analysis.

Supplement 15. differential gene expression analysis overview.

Supplement 16-22. differential gene expression and pathway analysis.

Supplement 23. most commonly mutated genes in dataset.

Additional file 3: Supplementary note.

Acknowledgements

Not applicable

Author's contributions

SB is the sole author of the manuscript. The author(s) read and approved the final manuscript.

Funding

Supported by Else-Kröner-Fresenius-Stiftung, Grant No. 2016-Kolleg-19, Deutsche Forschungsgemeinschaft (DFG), Grant No. EN 179/13-1 and Frauke Weiskam + Christel Ruranski-Stiftung, Grant No. T 0136 – 33.661. Open Access funding enabled and organized by Projekt DEAL.

Availability of data and materials

All raw data generated in this study is included in the supplement. Original sequencing datasets are available here: <https://dcc.icgc.org/>. To access original sequencing data, a request has been made to the ICGC DACO, which will be granted according to the conditions set out here: <https://icgc.org/index.php?q=daco>.

Declarations

Ethics approval and consent to participate

For each study which was included in this reanalysis of data, appropriate ethics approval has been granted and all patients consented to participate. Detailed information is available here: <https://dcc.icgc.org/>.

Consent for publication

Not applicable

Competing interests

The author declares that he has no competing interests

Author details

¹Department I of Internal Medicine, Center for Integrated Oncology Aachen Bonn Cologne Duesseldorf, University of Cologne, Cologne, Germany.

²Cancer Center Cologne Essen – Partner Site Cologne, CIO Cologne, University of Cologne, Cologne, Germany. ³German Hodgkin Study Group, Cologne, Germany.

Received: 29 April 2020 Accepted: 18 February 2021

Published online: 22 April 2021

References

- Schwabe RF, Jobin C. The microbiome and cancer. *Nat Rev Cancer*. 2013; 13(11):800–12. <https://doi.org/10.1038/nrc3610>.
- Goodman B, Gardner H. The microbiome and cancer. *J Pathol*. 2018;244(5): 667–76. <https://doi.org/10.1002/path.5047>.
- zur Hausen H. Papillomaviruses and cancer: from basic studies to clinical application. *Nat Rev Cancer*. 2002;2:342–50.
- 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*. 2013;4(1):2513. <https://doi.org/10.1038/ncomms3513>.
- Cantalupo PG, Katz JP, Pipas JM. Viral sequences in human cancer. *Virology*. 2018;513:208–16. <https://doi.org/10.1016/j.virol.2017.10.017>.
- Niedobitek G, Pitteroff S, Herbst H, Shepherd P, Finn T, Anagnostopoulos I, et al. Detection of human papillomavirus type 16 DNA in carcinomas of the palatine tonsil. *J Clin Pathol*. 1990;43(11):918–21. <https://doi.org/10.1136/jcp.43.11.918>.
- Gillison ML, Koch WM, Capone RB, Spafford M, Westra WH, Wu L, et al. Evidence for a causal association between human papillomavirus and a subset of head and neck cancers. *J Natl Cancer Inst*. 2000;92(9):709–20. <https://doi.org/10.1093/jnci/92.9.709>.
- Walboomers JMM, Jacobs MV, Manos MM, Bosch FX, Kummer JA, Shah KV, et al. Human papillomavirus is a necessary cause of invasive cervical cancer worldwide. *J Pathol*. 1999;189(1):12–9. [https://doi.org/10.1002/\(SICI\)1096-9896\(199909\)189:1<12::AID-PATH431>3.0.CO;2-F](https://doi.org/10.1002/(SICI)1096-9896(199909)189:1<12::AID-PATH431>3.0.CO;2-F).
- Schwarz E, Freese UK, Gissmann L, Mayer W, Roggenbuck B, Stremlau A, et al. Structure and transcription of human papillomavirus sequences in cervical carcinoma cells. *Nature*. 1985;314(6006):111–4. <https://doi.org/10.1038/8314111a0>.
- Perz JF, Armstrong GL, Farrington LA, Hutin YJF, Bell BP. The contributions of hepatitis B virus and hepatitis C virus infections to cirrhosis and primary liver cancer worldwide. *J Hepatol*. 2006;45(4):529–38. <https://doi.org/10.1016/j.jhep.2006.05.013>.
- Shafritz DA, Shouval D, Sherman HI, Hadziyannis SJ, Kew MC. Integration of hepatitis B virus DNA into the genome of liver cells in chronic liver disease and hepatocellular carcinoma. *N Engl J Med*. 1981;305(18):1067–73. <https://doi.org/10.1056/NEJM198110293051807>.
- zur Hausen H. The search for infectious causes of human cancers: where and why. *Virology*. 2009;392:1–10.
- Uemura N, Okamoto S, Yamamoto S, Matsumura N, Yamaguchi S, Yamakido M, et al. Helicobacter pylori infection and the development of gastric cancer. *N Engl J Med*. 2001;345(11):784–9. <https://doi.org/10.1056/NEJMoa001999>.
- Watanabe T, Tada M, Nagai H, Sasaki S, Nakao M. Helicobacter pylori infection induces gastric cancer in mongolian gerbils. *Gastroenterology*. 1998;115(3):642–8. [https://doi.org/10.1016/S0016-5085\(98\)70143-X](https://doi.org/10.1016/S0016-5085(98)70143-X).
- Castellarin M, Warren RL, Freeman JD, Dreolini L, Krzywinski M, Strauss J, et al. Fusobacterium nucleatum infection is prevalent in human colorectal carcinoma. *Genome Res*. 2012;22(2):299–306. <https://doi.org/10.1101/gr.126516.111>.
- Kostic AD, Gevers D, Pedamallu CS, Michaud M, Duke F, Earl AM, et al. Genomic analysis identifies association of Fusobacterium with colorectal carcinoma. *Genome Res*. 2012;22(2):292–8. <https://doi.org/10.1101/gr.126573.111>.
- Repass J, Reproducibility Project: Cancer Biology, Iorns E, Denis A, Williams SR, Perfito N, et al. Replication Study: Fusobacterium nucleatum infection is prevalent in human colorectal carcinoma. *Elife*. 2018;7. <https://doi.org/10.7554/eLife.25801>.
- Hanahan D, Weinberg RA. Hallmarks of cancer: the next generation. *Cell*. 2011;144(5):646–74. <https://doi.org/10.1016/j.cell.2011.02.013>.
- Goodwin S, McPherson JD, McCombie WR. Coming of age: ten years of next-generation sequencing technologies. *Nat Rev Genet*. 2016;17(6):333–51. <https://doi.org/10.1038/nrg.2016.49>.
- Riley DR, Sieber KB, Robinson KM, White JR, Ganesan A, Nourbakhsh S, et al. Bacteria-human somatic cell lateral gene transfer is enriched in cancer samples. *PLoS Comput Biol*. 2013;9(6):e1003107. <https://doi.org/10.1371/journal.pcbi.1003107>.
- Khouri JD, Tannir NM, Williams MD, Chen Y, Yao H, Zhang J, et al. Landscape of DNA virus associations across human malignant cancers: analysis of 3,775 cases using RNA-Seq. *J Virol*. 2013;87(16):8916–26. <https://doi.org/10.1128/JVI.00340-13>.
- International Cancer Genome Consortium, T. I. C. G. et al. International network of cancer genome projects. *Nature* 464, 993–998 (2010).
- Gibbs RA, et al. A global reference for human genetic variation. *Nature*. 2015;526:68–74.
- Wood DE, Salzberg SL. Kraken: ultrafast metagenomic sequence classification using exact alignments. *Genome Biol*. 2014;15(3):R46. <https://doi.org/10.1186/gb-2014-15-3-r46>.
- McIntyre ABR, Ounit R, Afshinnekoo E, Prill RJ, Hénaff E, Alexander N, et al. Comprehensive benchmarking and ensemble approaches for metagenomic classifiers. *Genome Biol*. 2017;18(1):182. <https://doi.org/10.1186/s13059-017-1299-7>.
- Sczyrba A, Hofmann P, Belmann P, Koslicki D, Janssen S, Dröge J, et al. Critical Assessment of Metagenome Interpretation - a benchmark of metagenomics software. *Nat Methods*. 2017;14(11):1063–71. <https://doi.org/10.1038/nmeth.4458>.
- Caygill CP, Hill MJ, Braddick M, Sharp JC. Cancer mortality in chronic typhoid and paratyphoid carriers. *Lancet (London, England)*. 1994;343:83–4.
- Scanu T, Spaapen RM, Bakker JM, Pratap CB, Wu LE, Hofland I, et al. Salmonella manipulation of host signaling pathways provokes cellular transformation associated with gallbladder carcinoma. *Cell Host Microbe*. 2015;17(6):763–74. <https://doi.org/10.1016/j.chom.2015.05.002>.
- Geller LT, et al. Potential role of intratumor bacteria in mediating tumor resistance to the chemotherapeutic drug gemcitabine. *Science (80-)*. 2017; 357:1156–60.
- Feng Y, Ramnarine VR, Bell R, Volik S, Davicioni E, Hayes VM, et al. Metagenomic and metatranscriptomic analysis of human prostate microbiota from patients with prostate cancer. *BMC Genomics*. 2019;20(1): 146. <https://doi.org/10.1186/s12864-019-5457-z>.
- Kripalani-Joshi S, Law HY. Identification of integrated Epstein-Barr virus in nasopharyngeal carcinoma using pulse field gel electrophoresis. *Int J Cancer*. 1994;56(2):187–92. <https://doi.org/10.1002/ijc.2910560207>.
- Morton C, et al. Mapping of the human Blym-1 transforming gene activated in Burkitt lymphomas to chromosome 1. *Science (80-)*. 1984;223:173–5.
- Robinson KM, Crabtree J, Mattick JSA, Anderson KE, Dunning Hotopp JC. Distinguishing potential bacteria-tumor associations from contamination in a secondary data analysis of public cancer genome sequence data. *Microbiome*. 2017;5(1):9. <https://doi.org/10.1186/s40168-016-0224-8>.
- Thompson KJ, Ingle JN, Tang X, Chia N, Jeraldo PR, Walther-Antonio MR, et al. A comprehensive analysis of breast cancer microbiota and host gene expression. *PLoS One*. 2017;12(11):e0188873. <https://doi.org/10.1371/journal.pone.0188873>.
- Mazzoni E, Benassi MS, Corallini A, Barbanti-Brodano G, Taronna A, Picci P, Guerra G, D'Agostino A, Trevisiol L, Nocini PF, Casali MV, Barbanti-Brodano G, Martini F, Tognon M. Significant association between human osteosarcoma and simian virus 40. *Cancer*. 2015;121(5):708–15. <https://doi.org/10.1002/cncr.29137>. Epub 2014 Nov 6. PMID: 25377935.
- Ramanan P, Deziel PJ, Buschack NL. Gordonia bacteremia. *J Clin Microbiol*. 2013;51(10):3443–7. <https://doi.org/10.1128/JCM.01449-13>.
- Ding X, Yu Y, Chen M, Wang C, Kang Y, Li H, et al. Bacteremia due to Gordonia polyisoprenivorans: case report and review of literature. *BMC Infect Dis*. 2017;17(1):419. <https://doi.org/10.1186/s12879-017-2523-5>.
- Gupta M, Prasad D, Khara HS, Alcidi D. A rubber-degrading organism growing from a human body. *Int J Infect Dis*. 2010;14(1):e75–6. <https://doi.org/10.1016/j.ijid.2009.03.006>.
- Zenz T, Eichhorst B, Busch R, Denzel T, Häbe S, Winkler D, et al. TP53 mutation and survival in chronic lymphocytic leukemia. *J Clin Oncol*. 2010; 28(29):4473–9. <https://doi.org/10.1200/JCO.2009.27.8762>.
- Austen B, Powell JE, Alvi A, Edwards I, Hooper L, Starczynski J, et al. Mutations in the ATM gene lead to impaired overall and treatment-free survival that is independent of IGVH mutation status in patients with B-CLL. *Blood*. 2005;106(9):3175–82. <https://doi.org/10.1182/blood-2004-11-4516>.
- Balthazar EJ, Megibow AJ, Hulnick DH. Cytomegalovirus esophagitis and gastritis in AIDS. *AJR Am J Roentgenol*. 1985;144(6):1201–4. <https://doi.org/10.2214/ajr.144.6.1201>.
- Lepilier Q, Tripathy MK, Di Martino V, Kantelip B, Herbein G. Increased HCMV seroprevalence in patients with hepatocellular carcinoma. *Virol J*. 2011;8(1): 485. <https://doi.org/10.1186/1743-422X-8-485>.

43. Leonardsson H, Hreinsson JP, Löve A, Björnsson ES. Hepatitis due to Epstein–Barr virus and cytomegalovirus: clinical features and outcomes. *Scand J Gastroenterol.* 2017;52(8):893–7. <https://doi.org/10.1080/00365521.2017.1319972>.
44. Bruix J, Llovet JM. Hepatitis B virus and hepatocellular carcinoma. *J Hepatol.* 2003;39:59–63. [https://doi.org/10.1016/S0168-8278\(03\)00140-5](https://doi.org/10.1016/S0168-8278(03)00140-5).
45. Lai C-C, Cheng A, Liu WL, Tan CK, Huang YT, Chung KP, et al. Infections caused by unusual *Methylobacterium* species. *J Clin Microbiol.* 2011;49(9):3329–31. <https://doi.org/10.1128/JCM.01241-11>.
46. Langevin S, Vincelette J, Bekal S, Gaudreau C. First case of invasive human infection caused by *Cupriavidus metallidurans*. *J Clin Microbiol.* 2011;49(2):744–5. <https://doi.org/10.1128/JCM.01947-10>.
47. Ugge H, Udumyan R, Carlsson J, Andrén O, Montgomery S, Davidsson S, et al. Acne in late adolescence and risk of prostate cancer. *Int J Cancer.* 2018;142(8):1580–5. <https://doi.org/10.1002/ijc.31192>.
48. Davidsson S, Mölling P, Rider JR, Unemo M, Carlsson MG, Carlsson J, et al. Frequency and typing of *Propionibacterium* acnes in prostate tissue obtained from men with and without prostate cancer. *Infect Agent Cancer.* 2016;11(1):26. <https://doi.org/10.1186/s13027-016-0074-9>.
49. Cohen RJ, Shannon BA, McNeal JE, Shannon T, Garrett KL. *Propionibacterium acnes* associated with inflammation in radical prostatectomy specimens: a possible link to cancer evolution? *J Urol.* 2005;173(6):1969–74. <https://doi.org/10.1097/01.ju.0000158161.15277.78>.
50. Mahlen SD. *Serratia* infections: from military experiments to current practice. *Clin Microbiol Rev.* 2011;24(4):755–91. <https://doi.org/10.1128/CMR.00017-11>.
51. Fessler J, Matson V, Gajewski TF. Exploring the emerging role of the microbiome in cancer immunotherapy. *J Immunother Cancer.* 2019;7(1):108. <https://doi.org/10.1186/s40425-019-0574-4>.
52. Motzer RJ, Escudier B, McDermott D, George S, Hammers HJ, Srinivas S, et al. Nivolumab versus everolimus in advanced renal-cell carcinoma. *N Engl J Med.* 2015;373(19):1803–13. <https://doi.org/10.1056/NEJMoa1510665>.
53. Sung W-K, et al. Genome-wide survey of recurrent HBV integration in hepatocellular carcinoma. *Nat Genet Vol.* 2012;44:765–9.
54. Feng Q, Liang S, Jia H, Stadlmayr A, Tang L, Lan Z, et al. Gut microbiome development along the colorectal adenoma-carcinoma sequence. *Nat Commun.* 2015;6(1). <https://doi.org/10.1038/ncomms7528>.
55. Gopalakrishnan V, Helmink BA, Spencer CN, Reuben A, Wargo JA. The influence of the gut microbiome on cancer, immunity, and cancer immunotherapy. *Cancer Cell.* 2018;33(4):570–80. <https://doi.org/10.1016/j.ccell.2018.03.015>.
56. Luo C, Knight R, Siljander H, Knip M, Xavier RJ, Gevers D. ConStrains identifies microbial strains in metagenomic datasets. *Nat Biotechnol.* 2015;33(10):1045–52. <https://doi.org/10.1038/nbt.3319>.
57. Bharti R, Grimm DG. Current challenges and best-practice protocols for microbiome analysis. *Brief Bioinform.* 2019;1–16.
58. Welcome | ICGC Data Portal. Available at: <https://dcc.icgc.org/>. (Accessed: 6th June 2019)
59. Index von /vol1/ftp/. Available at: <ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/>. (Accessed: 6th June 2019)
60. European Nucleotide Archive < EMBL-EBI. Available at: <https://www.ebi.ac.uk/ena>. (Accessed: 6th June 2019)
61. Afgan E, Baker D, Batut B, van den Beek M, Bouvier D, Čech M, et al. The Galaxy platform for accessible, reproducible and collaborative biomedical analyses: 2018 update. *Nucleic Acids Res.* 2018;46(W1):W537–44. <https://doi.org/10.1093/nar/gky379>.
62. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence Alignment/Map format and SAMtools. *Bioinformatics.* 2009;25(16):2078–9. <https://doi.org/10.1093/bioinformatics/btp352>.
63. Picard Tools - By Broad Institute. Available at: <http://broadinstitute.github.io/picard/>. (Accessed: 6th June 2019)
64. Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics.* 2010;26(6):841–2. <https://doi.org/10.1093/bioinformatics/btq033>.
65. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics.* 2014;30(15):2114–20. <https://doi.org/10.1093/bioinformatics/btu170>.
66. Blankenberg D, Gordon A, von Kuster G, Coraor N, Taylor J, Nekrutenko A, et al. Manipulation of FASTQ data with Galaxy. *Bioinformatics.* 2010;26(14):1783–5. <https://doi.org/10.1093/bioinformatics/btq281>.
67. Rognes T, Flouri T, Nichols B, Quince C, Mahé F. VSEARCH: a versatile open source tool for metagenomics. *PeerJ.* 2016;4:e2584. <https://doi.org/10.7717/peerj.2584>.
68. The R Foundation for Statistical Computing R version 3.3.2. (<https://www.r-project.org/>).
69. Escapa IF, et al. New insights into human nostril microbiome from the expanded Human Oral Microbiome Database (eHOMD): a Resource for the Microbiome of the Human Aerodigestive Tract. *mSystems.* 2018;3:e00187–18.
70. Dewhirst FE, Chen T, Izard J, Paster BJ, Tanner ACR, Yu WH, et al. The human oral microbiome. *J Bacteriol.* 2010;192(19):5002–17. <https://doi.org/10.1128/JB.00542-10>.
71. HOMD: Human Oral Microbiome Database. Available at: <http://www.homd.org/index.php?name=HOMD>. (Accessed: 6th June 2019)
72. de Goffau MC, Lager S, Salter SJ, Wagner J, Kronbichler A, Charnock-Jones DS, et al. Recognizing the reagent microbiome. *Nat Microbiol.* 2018;3(8):851–3. <https://doi.org/10.1038/s41564-018-0202-y>.
73. Salter SJ, Cox MJ, Turek EM, Calus ST, Cookson WO, Moffatt MF, et al. Reagent and laboratory contamination can critically impact sequence-based microbiome analyses. *BMC Biol.* 2014;12(1):87. <https://doi.org/10.1186/s12915-014-0087-z>.
74. Jervis-Bardy J, Leong LEX, Marri S, Smith RJ, Choo JM, Smith-Vaughan HC, et al. Deriving accurate microbiota profiles from human samples with low bacterial content through post-sequencing processing of Illumina MiSeq data. *Microbiome.* 2015;3(1):19. <https://doi.org/10.1186/s40168-015-0083-8>.
75. Leon LJ, et al. Enrichment of clinically relevant organisms in spontaneous preterm-delivered placentas and reagent contamination across all clinical groups in a large pregnancy cohort in the United Kingdom. *Appl Environ Microbiol.* 2018;84:e00483–18.
76. Laurence M, Hatzis C, Brash DE. Common contaminants in next-generation sequencing that hinder discovery of low-abundance microbes. *PLoS One.* 2014;9(5):e97876. <https://doi.org/10.1371/journal.pone.0097876>.
77. Glassing A, Dowd SE, Galandiuk S, Davis B, Chiodini RJ. Inherent bacterial DNA contamination of extraction and sequencing reagents may affect interpretation of microbiota in low bacterial biomass samples. *Gut Pathog.* 2016;8(1):24. <https://doi.org/10.1186/s13099-016-0103-7>.
78. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+: architecture and applications. *BMC Bioinformatics.* 2009;10(1):421. <https://doi.org/10.1186/1471-2105-10-421>.
79. Wally N, Schneider M, Thannesberger J, Kastner MT, Bakonyi T, Indik S, et al. Plasmid DNA contaminant in molecular reagents. *Sci Rep.* 2019;9(1):1652. <https://doi.org/10.1038/s41598-019-38733-1>.
80. Embedding projector - visualization of high-dimensional data. Available at: <https://projector.tensorflow.org/>. (Accessed: 6th June 2019)
81. Morpheus. Available at: <https://software.broadinstitute.org/morpheus/>. (Accessed: 6th June 2019)
82. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics.* 2009;25(14):1754–60. <https://doi.org/10.1093/bioinformatics/btp324>.
83. Sondka Z, Bamford S, Cole CG, Ward SA, Dunham I, Forbes SA. The COSMIC Cancer Gene Census: describing genetic dysfunction across all human cancers. *Nat Rev Cancer.* 2018;18(11):696–705. <https://doi.org/10.1038/s41568-018-0060-1>.
84. Warden CD, Yuan Y-C, Wu X. Optimal calculation of RNA-Seq fold-change values. *Int J Comput Bioinform In Silico Model.* 2013;2:285–92.
85. Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, et al. Gene ontology: tool for the unification of biology. *Gene Ontology Consortium Nat Genet.* 2000;25(1):25–9. <https://doi.org/10.1038/75556>.

Publisher's Note

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