

RESEARCH ARTICLE

Open Access

# Reference-free transcriptome signatures for prostate cancer prognosis



Ha T.N. Nguyen<sup>1</sup>, Haoliang Xue<sup>1</sup>, Virginie Firlej<sup>2</sup>, Yann Ponty<sup>3</sup>, Melina Gallopin<sup>1</sup> and Daniel Gautheret<sup>1\*</sup>

## Abstract

**Background:** RNA-seq data are increasingly used to derive prognostic signatures for cancer outcome prediction. A limitation of current predictors is their reliance on reference gene annotations, which amounts to ignoring large numbers of non-canonical RNAs produced in disease tissues. A recently introduced kind of transcriptome classifier operates entirely in a reference-free manner, relying on k-mers extracted from patient RNA-seq data.

**Methods:** In this paper, we set out to compare conventional and reference-free signatures in risk and relapse prediction of prostate cancer. To compare the two approaches as fairly as possible, we set up a common procedure that takes as input either a k-mer count matrix or a gene expression matrix, extracts a signature and evaluates this signature in an independent dataset.

**Results:** We find that both gene-based and k-mer based classifiers had similarly high performances for risk prediction and a markedly lower performance for relapse prediction. Interestingly, the reference-free signatures included a set of sequences mapping to novel lncRNAs or variable regions of cancer driver genes that were not part of gene-based signatures.

**Conclusions:** Reference-free classifiers are thus a promising strategy for the identification of novel prognostic RNA biomarkers.

**Keywords:** Reference-free transcriptomic, Supervised learning, Prostate cancer signature

## Introduction

The outcome of human cancer can be predicted in part through gene expression profiling [1–3]. Outcome prediction is particularly important in prostate cancer (PCa), where distinguishing indolent from aggressive tumors would prevent unnecessary treatment and improve patients' quality of life. However, currently there is no reliable signature of aggressive prostate cancer. Pathologists classify prostate tumor biopsies using scoring systems such as the Gleason score that evaluates tumor differentiation and Tumour, Node, Metastasis (TNM) staging that evaluates tumor extent and propagation. Gleason, TNM and Prostate-specific antigen (PSA) levels

can be combined into a low, medium or high risk status [4]. Several studies used gene expression profiles to derive predictors of Gleason score or risk [5–8]. Other studies predicted actual clinical progression (tumor recurrence or metastasis) after several years of patient followup. Clinical progression can be evaluated either indirectly through monitoring of PSA levels (BCR=biochemical relapse) [9–12] or upon direct clinical observation [13–16]. Gene expression predictors usually take the form of a signature, that is a set of genes or transcripts and associated coefficients of a model that can be used to predict risk or outcome from a patient sample. Commercial tests such as Decipher and Oncotype DX predict prostate cancer risk based on gene expression. However these are still not recommended for routine use [17]. In general, the prostate cancer community has progressed pretty well at identifying low and high risk patients, but men with mid-range

\*Correspondence: [daniel.gautheret@universite-paris-saclay.fr](mailto:daniel.gautheret@universite-paris-saclay.fr)

<sup>1</sup>Institute for Integrative Biology of the Cell, UMR 9198, CEA, CNRS, Université Paris-Saclay, Gif-Sur-Yvette, France

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.

risk face more uncertainty and would most benefit from improved tests.

Gene expression profiling of prostate biopsies is performed either using DNA microarrays [13–16] or high throughput RNA sequencing (RNA-seq) [5–8]. An important advantage of RNA-seq is its ability to identify novel genes or transcripts, which can in principle be incorporated into predictive signatures. However, RNA-seq analysis is usually performed in a “reference-based” fashion, ie. by using RNA-seq reads to quantify a predetermined set of transcripts. This amounts to using RNA-seq in the same way as a microarray that only quantifies a predetermined set of probes. Yet, there is abundant evidence that non-reference RNAs are frequent in disease tissues and may constitute clinically useful biomarkers [18]. Therefore one may expect that prognostic models incorporating non-reference RNAs may carry substantial benefits.

Our group [19, 20] and others [21] introduced new k-mer based strategies to analyse RNA-seq data in a “reference-free” manner, that is without mapping sequence reads to a predefined set of genes or transcripts. K-mers are sub-sequences of fixed length which are extracted and quantified from sequence files. When applied to medical RNA-seq datasets using appropriate statistical methods, this strategy identifies any sub-sequence whose increased abundance is associated to a given clinical label. This may include novel splice variants, long non-coding RNAs (lncRNAs) or RNAs from repeated retroelements [19, 20] which are ignored by conventional protocols based on reference gene annotations.

Although attractive in principle, k-mer derived prognostic signatures pose two major challenges. First, a single RNA-seq dataset commonly contains tens to hundreds of millions distinct k-mers. Therefore false positive and replicability issues encountered with gene expression profiles [22–25] are expected to worsen with k-mer count matrices. The second challenge is related to the transfer of a k-mer signature across independent datasets. Signatures inferred from an initial discovery set are expected to generalize to any independent dataset. In the absence of a unifying gene concept, independent validation requires matching signature k-mers to read sequences from the new dataset. This may cause significant signal loss if sequencing or library preparation technologies differ.

Our main objective here was to compare the characteristics and performances of reference-based and reference-free classifiers for PCa risk and relapse prediction. We built both types of classifiers using the same discovery dataset and assessed their performances in independent datasets using equivalent pipelines and parameters. For the reference-free approach, this required special developments to reduce the number of variables and to transfer expression measures between datasets. We present

below a detailed analysis of the relative performances and sequence contents of the different classifiers and discuss possible future developments to improve performances of models.

## Materials and methods

### Data acquisition and outcome labelling

We used tumor samples from the TCGA-PRAD data collection [26] (N=505) for signature discovery. The resulting classifiers were then assessed in two independent datasets, from the Canadian Prostate Cancer Genome Network (ICGC-PRAD-CA) [27] (N=148) and from the Portuguese Oncology Institute’s “Porto” cohort, analyzed in Stelloo et al. [28] (N=91). All three datasets were produced from radical prostatectomies and used similar technologies for library preparation (frozen samples, poly(A)+ RNA selection) and Illumina sequencing, however they differed by read-size, read depth, strandedness and use of single or paired ends sequencing (Table 1).

TCGA-PRAD RNA-seq data were retrieved from dbGAP accession phs000178.v9.p8 with permission. ICGC-PRAD-CA RNA-seq data (EGAD00001004424) were downloaded from the European Genome-Phenome Archive (EGA) with permission. The RNA-seq files from the “Porto” cohort [28] were retrieved from GEO, under accession GSE120741. Clinical information was retrieved from Liu et al. [29] for TCGA-PRAD, from Fraser et al. [27] for ICGC-PRAD and from sample metadata of GEO accession GSE120741 for Stelloo et al. [28].

We built predictors for risk and relapse using two-class prediction models. To achieve a clear separation between the two classes, we only focused on high risk (HR) samples versus low risk (LR) samples, ignoring the medium risk, and we focused on relapse prior to a given year and non-relapse after a given year. For this reason, only a fraction of samples could be labelled for a given class in each set. Risk information was not available in the Stelloo dataset and relapse labelling on the ICGC dataset led to a small validation set (only 7 relapse samples).

We classified tumor specimens into low-risk and high-risk groups using an adaptation of d’Amico’s classification which does not take into account the PSA rate but only the anatomo-pathological data on the basis of Gleason and TNM features as performed previously [20]. Tumors with Gleason score 6/7 (3+4) and TNM stage pT1/2 were classified as low risk. Tumors with Gleason score 8/9 and/or TNM stage pT3b/4 were defined as high-risk. Tumors classified as pT3a, pT1 or (pT2 and Gleason (4+3)) were considered as intermediate and excluded from the analysis. 374 TCGA-PRAD tumors and 63 ICGC-PRAD-CA tumors could be labelled for LR or HR. We could not obtain Gleason/TNM scores for Stelloo et al, hence we did not annotate risk for this cohort.

**Table 1** Characteristics of prostate tumor RNA-seq datasets

Study	RNA-seq library type	Reads/sample	#Tumor samples	Risk		Relapse	
				LR	HR	NO	YES
TCGA-PRAD	Poly(A)+ unstranded 2x50nt	130M	505	134	240	56	58
ICGC-PRAD	Poly(A)+ stranded 2x100nt	313M	148	40	23	49	7
STELLOO	Poly(A)+ stranded 1x65nt	20M	91			43	48

For relapse analysis, we distinguished patients with biochemical relapse (BCR) and time to BCR <2yr and patients with no BCR after 5 years or longer, except for Stelloo et al. where only precomputed relapse data was available with cutoffs at 5yr and 10yr, respectively (Table 2). BCR information was obtained from Table S1 of Liu et al. [29] for TCGA-PRAD and from table S1 (PFS field) of Fraser et al. [27] for ICGC-PRAD. Precomputed relapse data for Stelloo et al. was taken from SRA accession PRJNA494345.

**A generic framework to infer reference-based and reference-free signatures**

Risk and relapse predictors were derived using a combination of feature selection and supervised learning (Fig. 1). The predictive model was tuned over a discovery (or training) dataset and its performance was then evaluated on an independent validation (or testing) dataset, to avoid selection bias [30]. The same procedure was used for reference-based and reference-free models, however two extra steps were included to obtain and validate reference-free signatures. First a procedure was implemented to reduce the k-mer matrix using a sequence assembly-like algorithm to merge k-mers into contigs based on their sequence overlap and on the similarity of their count vectors. This step led to a contig count table an order of magnitude smaller than the initial k-mer count table (see “Results” section below). Feature selection and model fitting were performed over this contig table. A second adaptation was necessary to validate the reference-free signature in an independent dataset. This required extracting k-mers from both the signature and the sequence files of the independent set, and compute the signature expression in the independent set based on counts of matching k-mers. The pipeline is detailed in Methods. Note that we

select features and train a predictive model only on the discovery dataset. The model is then applied to the validation set with no retraining (i.e. with the same coefficients) for an unbiased evaluation of the signature.

**Gene and k-mer count matrices**

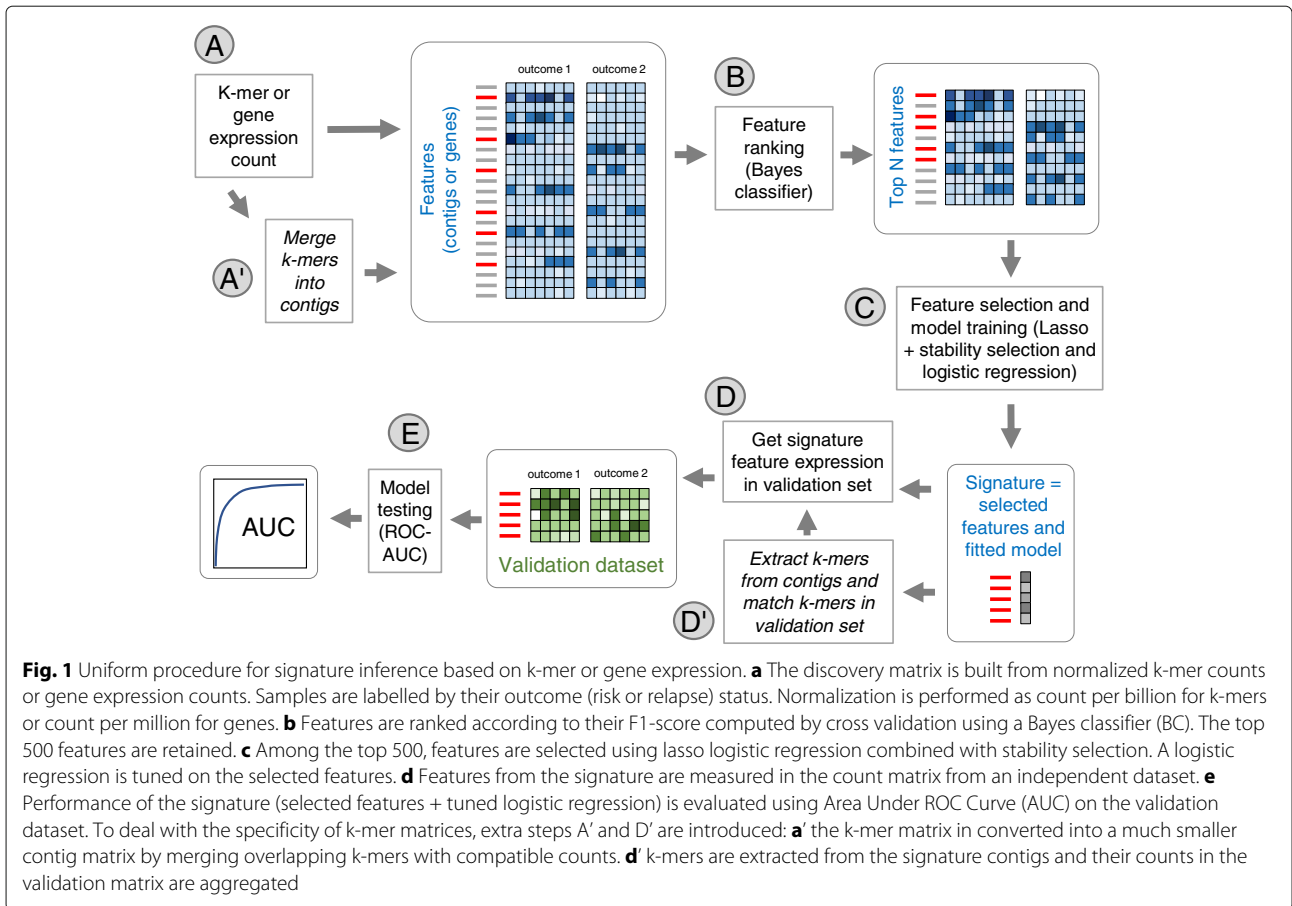
DEkupl-run [19] was used to produce gene and k-mer count matrices for each dataset. DEkupl-run converts FASTQ files to k-mer counts using Jellyfish [31], joins individual sample counts into a single count table and filters out low count k-mers. K-mer size was set to 31, lib\_type to unstranded, and parameters min\_recurrence and min\_recurrence\_abundance were set for each dataset as in Additional file 4: Table S1. K-mer size was set to 31 as commonly adopted for human transcriptome applications [19, 32]. Note that contrary to TCGA-PRAD, ICGC-PRAD uses stranded RNA-seq libraries. However we could not use this information as signatures were produced from unstranded libraries. We thus built all k-mer tables in canonical mode, which amounts to consider all libraries as unstranded. Gene expression was computed using Kallisto v0.43.0 [32] with Gencode V24 as a reference transcriptome. Gene-level counts were obtained by summing counts for all transcripts of each gene. Gene expression matrices were submitted to the same recurrence filters as k-mer tables to remove low expression genes. After count tables were generated and filtered, the k-mer merging and differential expression analysis module of DEkupl-run were not used. Instead, tables were further processed as explained below.

**Reduction of k-mer matrix via contig extension**

k-mer occurrence tables were converted into contig occurrence tables using an extension procedure similar

**Table 2** Relapse group definitions

Relapse group	TCGA-PRAD	ICGC-PRAD	STELLOO
Relapse (YES)	PFS = 1 and	BCR = “Yes” and	BCR = “Yes” and
	PFS.time <2yr	BCR.time <2yr	BCR.time <5yr
Non relapse (NO)	PFS = 0 and	BCR = “No” and	BCR = “No” and
	PFS.time >5yr	BCR.time >5yr	BCR.time >10yr



to that described in Audoux et al. [19]. We define here as contig any sequence produced by merging 1 or more k-mers. Briefly, contigs overlapping by (k-1) to (k-15) nucleotide were iteratively merged into longer contigs till any of the following condition was encountered. In a straightforward case, extension stops when no more overlapping contig is available. Alternatively, extension stops when ambiguity is introduced i.e. when competing extension paths occur. Lastly, we applied here an intervention not included in Audoux et al. [19] by considering sample count compatibility between contigs, as shown in Fig. 2. Sample count compatibility is measured by the mean value of absolute contrast (MAC) between the counts of the two contigs across all samples, i.e.

$$MAC(c_1, c_2) = \text{mean}_{s \in \{\text{samples}\}} \left( \left| \frac{c_{1,s} - c_{2,s}}{c_{1,s} + c_{2,s}} \right| \right)$$

where  $c_1$  and  $c_2$  are count vectors of two contigs to be merged, and  $c_{1,s}$  and  $c_{2,s}$  are counts in sample  $s$  from the corresponding count vectors. The extension is rejected if  $MAC > 0.25$ . In this way, all contigs are guaranteed to have member k-mers with consistent sample count vec-

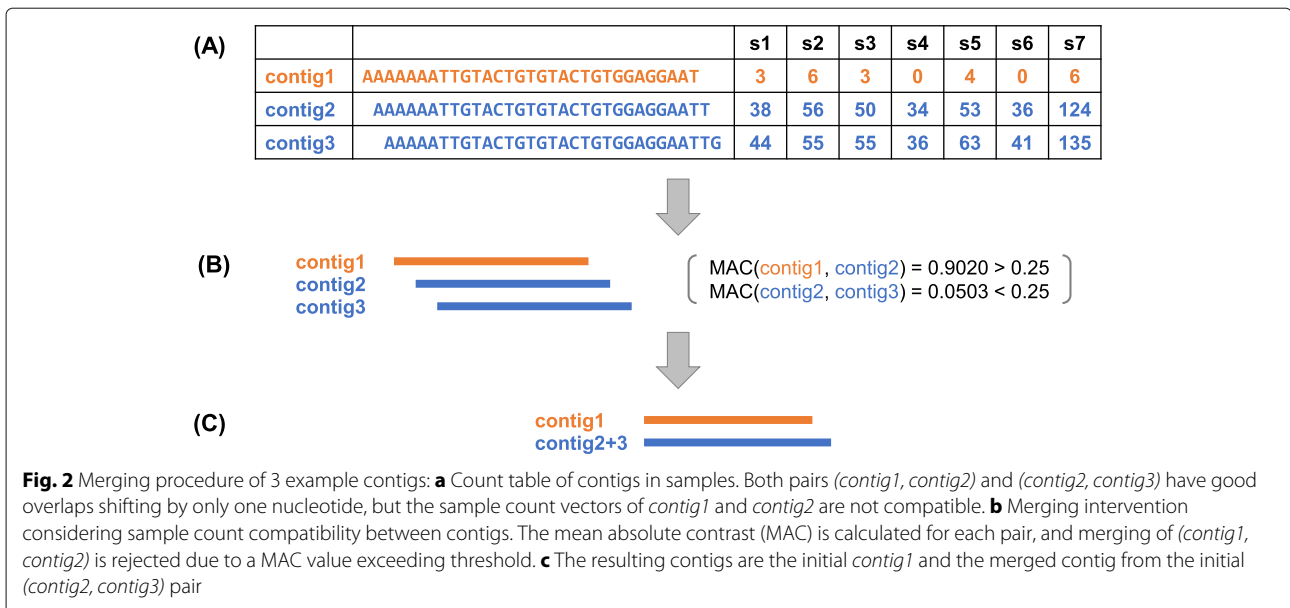
tors. After the merging procedure, the new contig's sample count vector is set to the mean of composite k-mer's sample count vectors.

**Count normalization**

To account for differences in sequencing depth among samples, we applied a normalization step on feature counts (genes or contigs) in discovery and validation datasets. Each feature count in a sample is divided by the sum of all feature counts in this sample, then multiplied by a constant base number:

$$e_{f,s} \leftarrow \frac{e_{f,s}}{\sum_{f \in \{\text{features}\}} e_{f,s}} \cdot C_b,$$

where  $e_{f,s}$  refers to count of feature  $f$  in sample  $s$ , and  $C_b$  is the base constant. For genes,  $C_b = 10^6$  resulting in a conventional count per million (CPM) normalization, while for contigs, we used  $C_b = 10^9$ , or count per billion (CPB). For contigs, normalization is applied on the contig count table produced after contig extension and for genes it is applied on the recurrence filtered gene expression matrix.



### Univariate features ranking

Given the limited number of samples, it was necessary to reduce the number of features (genes or contigs) in the dataset. We discarded irrelevant features to focus on a subset of 500 top candidates for subsequent feature selection. To rank features, we selected a Bayes classifier because the C++ implementation of this classifier was the fastest to run among several available feature ranking tools. We did not try to optimize this part to avoid biasing the comparison towards gene-based or gene-free methods. In detail, we performed prediction of status (risk/relapse) using a Bayes classifier on each independent feature, after log transformation of the normalized counts (after adding an offset 1 to avoid numerical problem). To assess the quality of the prediction, we computed the average  $f_1$  score by 5-fold cross validation ( $f_1 = \frac{2 \cdot \text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$ , where  $\text{precision} = TP / (TP + FP)$  and  $\text{recall} = TP / (TP + FN)$  and  $FP, TP, FN$  are respectively the False Positive, True Positive and False Negative). In cases where 5-fold cross-validation returned an undefined value,  $f_1$  score was set to 0 (the worst). The average  $f_1$  score was used to rank features. The Bayes classifier implementation was taken from the MLPack library [33].

### Feature selection, model fitting and predictor evaluation

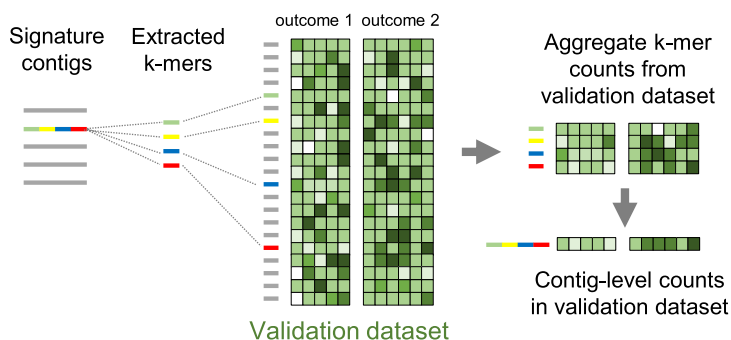
To select a subset of non-correlated features (genes or contigs) among the top 500 candidates, we performed penalized logistic regression using the implementation from the glmnet R package [34]. We implemented stability selection [35]: only features selected with a frequency of being selected above 0.5 upon 2000 resamples of the input dataset were retained. To evaluate the performance

of the selected features on the discovery (training dataset), we fitted a logistic regression and computed the area under the ROC curve (AUC) using a 10-fold cross validation scheme, repeated 20 times, as implemented in the caret package [36]. To handle imbalanced datasets, we included optional oversampling and downsampling in our evaluation procedures [37]. We also computed the Precision-Recall AUC, a more informative metric than the ROC AUC when evaluating binary classifiers on imbalanced datasets [38]. To assess the performance of the signature on the external validation datasets, we fitted a logistic regression on the whole discovery dataset and applied the predictor to the validation datasets. In the reference-free approach, some features present in the signature were not found in the validation (see below). In this case, the coefficient of the logistic regression corresponding to missing features were set to zero. Signature contigs were annotated through BLAST alignment vs. Gencode V34 transcripts. HGNC symbols for signature genes were obtained from the Ensembl EnsDb.Hsapiens.v79 R package [39].

### Matching signature contigs in the validation cohort

To measure contig expression in the validation cohort we implemented the procedure schematized in Fig. 3. The procedure comprises two main steps: (1) all k-mers from signature contigs were extracted and identified in the k-mer count matrix generated from the validation cohort and (2) the resulting sub-matrix was used to estimate each contig's expression in the validation cohort, measured for each sample as the median of extracted k-mer counts.





**Fig. 3** Procedure for inferring signature contig expression in an independent validation dataset. The colored contig from the signature is quantified in the validation cohort by extracting all its constituent k-mers and retrieving the corresponding k-mer counts from validation k-mer count matrix. The count vector of the contig in each sample of the validation dataset is taken as the median of counts for k-mers in this sample

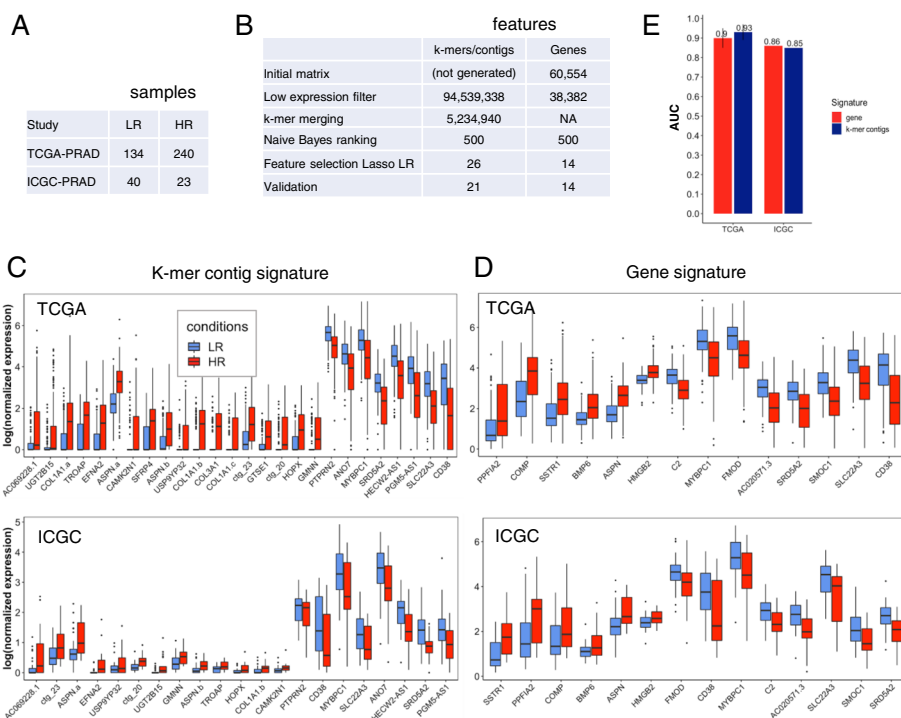
**Results**

**A reference-free risk signature for prostate cancer**

We first applied the gene-free and gene-based signature discovery procedures detailed above to infer PCa risk signatures. The k-mer table for 374 TCGA-PRAD risk-labelled samples (Fig. 4a) had 94M k-mers after low count filtering. The merging step reduced it to 5.2M contigs, i.e. achieving a considerable 18-fold reduction in size (Fig. 4b).

Contig sizes (mean=49nt, median=34nt, Table 3) were small relatively to a typical human RNA, which is characteristic of the adopted contig extension procedure [19] (see “Reduction of k-mer matrix via contig extension” section).

The 5.2M contig matrix and the 38k gene expression matrix were submitted to screening using univariate Bayes classification and the top scoring 500 features were retained for feature selection and model fitting. Inter-



**Fig. 4** Risk signatures generation and analysis. **a** Characteristics of prostate tumor RNA-seq datasets. **b** Result of filtering procedure on the k-mer and gene matrices for risk analysis. Expression of risk signature elements in LR and HR samples in the TCGA-PRAD and ICGC-PRAD cohorts **c** k-mer contig signature; **d** Gene signature. **e** Signature performances for risk prediction in the TCGA-PRAD and ICGC-PRAD cohorts

**Table 3** Contig sizes (Risk model)

	After k-mer merging	After Bayes classifier ranking
Mean contig size (nt)	49.1	189
Median contig size (nt)	34	61

tingly, the 500 top scoring contigs were significantly longer than prior to selection (median 61nt vs. 34nt, Table 3), suggesting the procedure tended to eliminate spurious short contigs.

Finally, Lasso logistic regression produced a reference-free signature of 26 contigs and a reference-based signature of 14 genes (Fig. 4b). Ten-fold cross validation performances of both signatures were very high on the discovery dataset (0.90 and 0.93 for genes and k-mers, respectively) (Fig. 4e), which is an over-estimated performance since features here were tested on the same dataset used to select features [30]. PR-AUC and ROC-AUC on different sampling techniques to adjust the class distribution of a dataset are also presented in Additional file 4: Table S2. These results lead to the same conclusion as the ones presented in (Fig. 4e).

Figure 4c shows the 26 contigs in the reference-free risk signature and their abundance distribution in LR and HR samples. 24/26 contigs mapped Gencode transcripts from 21 unique genes (Additional file 1). Eleven of the 21 genes were also found in a list 180 genes compiled from published PCa outcome signatures (Additional file 2), which is a highly significant enrichment ( $P$ -value =  $7.9e-9$ , Fisher’s exact test), especially when considering that no gene information was used to infer our signature. The gene and contig signatures involved five shared genes: MYBPC1, ASPN, SLC22A3, SRD5A2 and CD38 (Additional file 2, Fig. 4c and d). The first four genes are part of published prostate risk signatures. CD38 is particular in that it is the most downregulated in both signatures and it is not part of previous signatures. However, downregulation of this gene has been associated with poor outcome in prostate cancer [40], supporting its status as a high risk biomarker. Risk signature contigs mapped at least five other genes with established driver roles in PCa or other cancers: CAMK2N1 [41], COL1A1 [42], GTSE1 [43] and PTPRN2 [44], supporting the relevance of these sequence contigs in PCa etiology.

Of the two contigs that did not map any Gencode transcript, one aligned to an intron of GMNN (ctg\_20), a gene also mapped by an exonic contig, the other an intron of LDLRAD4 (ctg\_23). Contig ctg\_23 corresponds to a 1.29 kb spliced transcript located between exons 4 and 5 of LDLRAD4 and is strongly upregulated in HR samples, as displayed in the Integrative Genomics Viewer (IGV) [45] in Additional file 4: Figure S1. Although

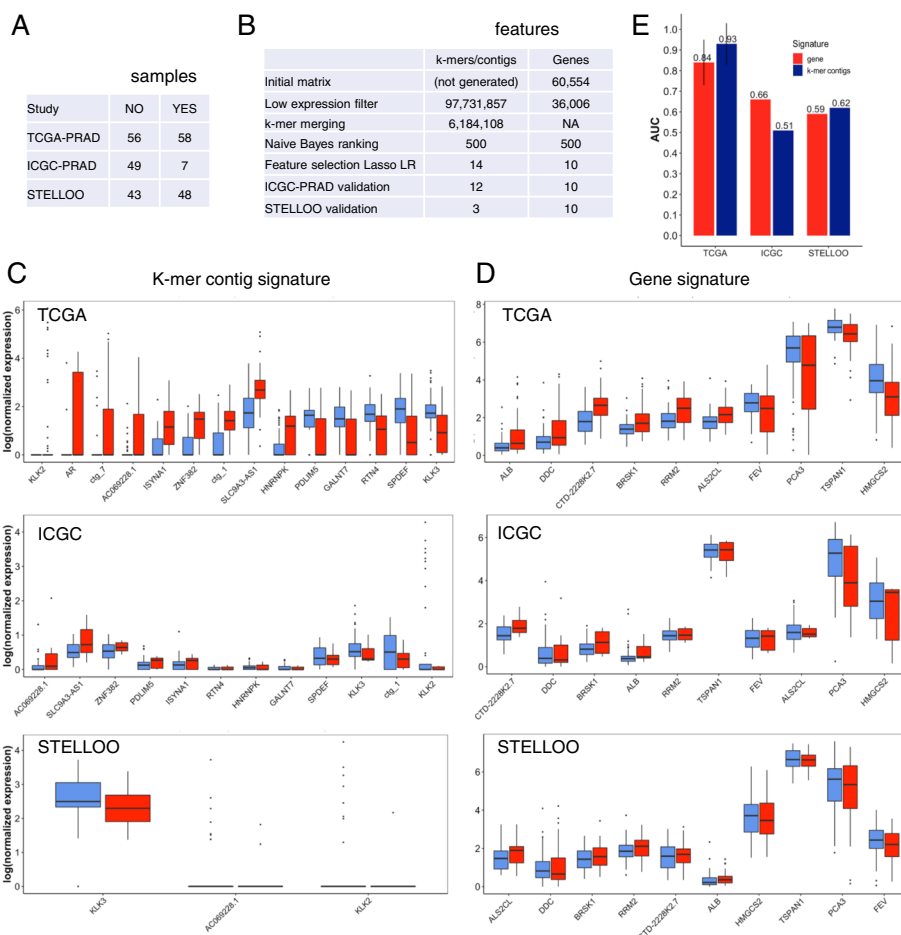
ctg\_23 partly maps short annotated LDLRAD4 isoforms, its expression seems unrelated to that of the longer LDLRAD4 transcripts whose coverage in flanking exons is 4-6 times lower than ctg\_23 (Additional file 4: Figure S2.) Therefore ctg\_23 likely comes from an independent lncRNA. The host gene LDLRAD4 is a negative regulator of TGF-beta signaling with roles in proliferation and apoptosis and was recently associated to negative outcome in other tumor types [46, 47]. Lastly, one contig (ctg\_11, EFNA2) was probably misassigned to the EFNA2 gene since it maps to a highly expressed discrete area just 3’ of EFNA2 while EFNA2 seems silent. Thus ctg\_11 probably comes from an independent lncRNA as well (Additional file 4: Figure S3).

To assess the replicability of risk signatures, we evaluated their performance in the ICGC-PRAD independent dataset. To this aim, we developed a specific procedure to estimate the expression of an arbitrary sequence contig across datasets using matched k-mers (see “Materials and methods” section). The 26 contigs represented 1444 k-mers, of which 97% were present in the ICGC-PRAD validation dataset. Overall 5 contigs (SFRP4, GTSE1, COL3A1, COL1A1.a, COL1A1.c) could not be quantified in the validation set due to lack of supporting k-mers (see Fig. 4b and c). In spite of this, the reference-free signature had similar performance in the validation set as the reference-based signature (0.85 and 0.86 respectively, Fig. 4e), although the later did not sustain any loss when transferred to the independent cohort (Fig. 4b). High prediction AUCs observed in the independent validation cohorts indicate a strong replicability of both the reference-free and reference-based risk signatures.

**Relapse signatures contain key PCa drivers**

For relapse prediction, we distinguished patients with biochemical relapse within less than 2 years and patients with no BCR after 5 years or longer. Application of the gene-free and gene-based signature discovery procedures to relapse prediction produced a 14-contig reference-free signature and a 10-gene reference-based signature (Additional file 2, Fig. 5b, c and d). The reference-free signature was populated by obvious PCa drivers. Strikingly, 3 contigs matched KLK2, AR and KLK3, which are among the most important genes in PCa onset and progression [48], the androgen receptor (AR) and two of its main targets, KLK2 and KLK3, the later encoding the PSA protein (Fig. 5c). Another contig matched SPDEF, a gene whose loss is associated to PCa metastasis [49].

Contigs matching KLK2 and AR were overexpressed 23-fold and 7-fold, respectively in relapsed patients while the contig matching KLK3 was depleted 1.8 fold. The AR contig matches exon 1 of AR and contains an non-templated poly-A end but no visible polyadenylation signal. The



**Fig. 5** Relapse signatures generation and analysis. **a** Characteristics of prostate tumor RNA-seq datasets. **b** Result of filtering procedure on the k-mer and gene matrices for relapse analysis. Expression of relapse signature elements in LR and HR samples in the TCGA-PRAD, ICGC-PRAD and STELLOO cohorts **c** k-mer contig signature; **d** Gene signature. **e** Signature performances for relapse prediction in the TCGA-PRAD, ICGC-PRAD and STELLOO cohorts

KLK2 contig is intronic and harbours a common SNP (rs62113074). The KLK3 contig is located in a distal part of the 3' UTR region present only in longer isoforms of KLK3. Its lower expression in relapsed patients was unexpected as low expression of PSA is usually associated to a lower risk. It is possible though that only this longer isoform is depleted in relapsing samples. The expression boxplot shows the KLK2 contig occurs only in a few outlier patients while the AR and KLK3 contigs are common (Fig. 5c). The contig matching SPDEF is a special variant of the 3' exon including two nonsynonymous SNPs. The SPDEF gene as a whole was highly expressed in both relapse and non-relapse samples but the contig expression was twice lower in average in relapse samples. Two contigs matched no known transcript: ctg\_7 is a low complexity sequence of unknown origin and ctg\_1 matches an intron of RPL9.

The contig matching lncRNA AC069228.1 also raised our attention since AC069228.1 is the only gene mapped by contigs in both relapse and risk signatures. The AC069228.1 lncRNA is antisense of PPFIA2, a protein tyrosine phosphatase that is itself an alleged urine biomarker of PCa [50]. The contigs from risk and relapse models match different regions of AC069228.1 (Figure S4). One is spliced, the other is a continuous 864 bp segment of a long exon. In both cases, a negative outcome (HR or relapse) is associated to a clearly higher expression of the contig, while the antisense gene PPFIA2 does not appear to follow the same trend (Figure S4).

Of note, the 10 genes in the reference-based signature were also clearly PCa-related: one was the major PCa biomarker PCA3 [51] and 5 others (DDC, RRM2, FEV, TSPAN1, HMGCS2) are involved in PCa etiology [52–56]. Therefore both gene-based and gene-free relapse signa-



tures were significant in terms of PCa related functions of their component genes or contigs.

### Relapse signatures do not accurately classify independent cohorts

Contrary to the risk signatures, relapse signatures showed little overlap with each other and with published PCa signatures (Additional file 2). Only PCA3 and KLK2 were found in prior signatures [16, 57] and the only gene found shared between relapse and risk signatures in this study was AC069228.1. The poor overlap in this study was not unexpected as the discovery samples for risk and relapse information were quite disjointed and not always consistent: for instance only 25% of the high risk samples were labelled for relapse and 28% of these did not relapse. Conversely, 51% of non-relapse patients were labelled as HR. Therefore risk and relapse classifiers were trained to recognize quite different phenotypes.

As in the risk model, both reference-based and reference-free signatures had excellent cross-validation performance on the discovery set (AUC of 0.84 and 0.93 respectively, Fig. 5e). However this should again be considered as an overly optimistic estimation due to the experimental design. Indeed, performances of both relapse signatures on the ICGC-PRAD and Stelloo validation sets were much lower (AUC 0.51 to 0.66), bordering randomness and confirming overfitting of the trained signatures. Substituting the logistic Regression classifier by Random Forest, or Boosted Logistic Regression did not improve performance of either model (Table S3). The reference-based model performed slightly better over ICGC-PRAD, and the reference-free model was slightly better over the Stelloo dataset (Fig. 5e). Furthermore, several genes and contigs in the discovery signatures had inconsistent expression variations in the validation datasets (Fig. 5c and d, Additional file 3). Overall two genes from the reference-based signature (ALB and CTD-2228K2.7) and 5 contigs from the reference-free signature (KLK2, AC069228.1, PDLIM5, RTN4, ctg\_1) changed logFC sign between the discovery and either validation cohort. This problem, which was not observed in risk models, underlines the poor replicability of the relapse signatures, whether or not reference-free.

Low replicability of the relapse model may be caused in part by weaknesses in validation datasets: the ICGC dataset had only 7 samples labelled for relapse (Fig. 5a) and the Stelloo dataset had very low coverage (Fig. 5a) which caused considerable loss when computing contig expression. Only three of the 14 signature contigs (AC069228.1, KLK2 and KLK3) could be quantified in the Stelloo dataset (Fig. 5b and c). Yet, we note that in spite of this loss the reference-free model still outperformed the reference-based model on this set (AUC of 0.62 vs.

0.59, Fig. 5e). Other limitations of the relapse model are addressed in the discussion.

## Discussion

### Properties of reference-free signatures

We evaluated here a method for building transcriptome classifiers that are totally reference-free, i.e. that do not require prior knowledge of genes or genome. The major interest of this approach lies in its ability to discover and incorporate in models previously unknown RNA biomarkers. Multiple examples exist of such disease-specific RNAs produced by genome alterations or deficient RNA processing and we hypothesized their inclusion in predictive models would be beneficial [18]. Applying a reference-free strategy to PCa outcome prediction, we obtained signatures made of short RNA contigs (median size 33 to 45 nt). These contigs are not full transcript models as can be produced by usual *de novo* assembly procedures. Instead, they often match SNPs or splice variants thus describing specific genetic or transcriptional events enriched in a patient group. Our strategy thus identifies RNA variations independently instead of lumping them into a full transcript model. Yet, the mapped genes were highly relevant to PCa etiology and included known cancer drivers LDLRAD4, GMNN, COL1A1, CD38, PTPRN2, GTSE1 and CAMK2N1 in the risk signature and KLK2, AR, KLK3, SPDEF in the relapse signature. Furthermore the risk signature comprised contigs matching two potential novel lncRNAs, located within LDLRAD4 and immediately downstream of EFNA2.

To our knowledge the only other software using a reference-free approach for inferring predictive signatures is Gecko [21]. Gecko uses machine learning (genetic algorithm) directly on the k-mer count matrix while we first reduce the matrix by grouping k-mers into contigs, before classification and machine learning. This enabled us to produce a signature composed of sequences larger than k, hence easier to interpret and quantify in an independent dataset.

Transferring a reference-free model to a new dataset is challenging. This requires that important features, such as SNPs, are precisely evaluated in the independent dataset. To this aim, we transferred signatures between datasets based on exact k-mer matches. As k-mer contents vary a lot between library preparation protocols, we expected this strategy to show poor sensitivity when discovery and validation datasets differed substantially. Indeed, transfer of signatures trained on the TCGA-PRAD dataset to the low coverage Stelloo dataset caused the loss of a majority of contigs. However, in this particular case, the remaining contigs were sufficient to maintain a prediction performance at the same level as that of the gene-based signature.

### Performances and generalization issues

To compare the reference-free and reference-based strategies, a common evaluation framework was adopted. For both risk and relapse predictions, performances of the reference-free classifiers were on a par with that of reference-based classifiers. However while risk signatures showed satisfying reproducibility, relapse signatures performed poorly in independent datasets.

A possible reason for the low performance of relapse models is our grouping of patients in discrete relapse and non relapse categories as done in other studies [9, 13, 15, 16]. This allowed us to address relapse prediction using the same logistic regression method as for risk, however this meant valuable patient information was left unused. A more accurate prediction of relapse may be achieved using survival models [10, 12, 14, 57, 58]. Adaptation of survival analysis tools to large k-mer matrices require additional developments that are certainly worth considering in the future.

A more general concern with relapse analysis is related to difficulty of predicting an outcome occurring several years after a sample is biopsied and analyzed. There might just be too little information available in the training data to infer a reliable classifier, a problem that is independent of the use of contigs or genes. However, both gene-level and contig-level signatures were highly enriched in PCa driver genes, which suggests information about tumor progression was indeed present in the primary tumor biopsy. The key problem with relapse analysis was more likely related to sample heterogeneity. The diversity of relapse mechanisms was not properly represented in a training set of 100 patients as we used here. Patient stratification have been proposed to deal with sample heterogeneity in omics data [59, 60]. Adaptations of these solutions to large k-mers matrices will also be considered in the future.

### Conclusion

For prediction of PCa risk and relapse, reference-free classifiers did not significantly outperform reference-based classifiers, however they incorporated a distinct set of RNA sequences including unannotated RNAs and novel variants of annotated RNAs. It is likely that with other diseases and datasets, novel biomarkers will be identified with an even greater impact on prediction performance. The reference-free approach will be of particular interest in problems where unknown RNAs are expected to play an important role, such as when studying rare diseases, poorly studied tissue types or when analysing dual human-pathogen RNA-seq samples. Our strategy also permits to infer efficient transcriptome classifiers in species lacking an accurate genome or transcriptome reference.

### Abbreviations

AUC: Are under the ROC curve; BCR: Biochemical relapse; HR: High risk; lncRNA: Long non-coding RNA; LR: Low risk; MAC: Mean absolute contrast; PCa: Prostate cancer; PSA: Prostate-specific antigen; RNA-seq: RNA sequencing; TNM: Tumour node metastasis

### Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s12885-021-08021-1>.

**Additional file 1:** Contig sequences and mapping locations in the risk and relapse signatures.

**Additional file 2:** Published PCa risk and relapse signatures. Genes in common between published and this publication's signatures.

**Additional file 3:** Contents and expression characteristics of all signatures in the discovery and validation datasets.

**Additional file 4:** Supplementary figures and tables.

### Acknowledgements

Not applicable.

### Authors' contributions

HTNN and HX developed the software, HTNN generated and analyzed the results, VF analyzed the clinical data, YP, MG and DG designed the experiments, MG and DG wrote the manuscript. All authors read and approved the final manuscript.

### Funding

This work was funded in part by Agence Nationale de la Recherche grant ANR-18-CE45-0020 and by a 911 Scholarship Fund from the Vietnamese Government to HTN.

### Availability of data and materials

The codes to reproduce the experiments are available on GitHub at: [https://github.com/i2bc/PCa-gene-based\\_vs\\_gene-free](https://github.com/i2bc/PCa-gene-based_vs_gene-free).

### Declarations

#### Ethics approval and consent to participate

Not applicable.

#### Consent for publication

Not applicable.

#### Competing interests

The authors declare that they have no competing interests.

#### Author details

<sup>1</sup>Institute for Integrative Biology of the Cell, UMR 9198, CEA, CNRS, Université Paris-Saclay, Gif-Sur-Yvette, France. <sup>2</sup>Institute of Biology, Université Paris Est Creteil, Creteil, France. <sup>3</sup>LIX CNRS UMR 7161, Ecole Polytechnique, Institut Polytechnique de Paris, Palaiseau, France.

Received: 18 December 2020 Accepted: 9 March 2021

Published online: 12 April 2021

### References

- Perou CM, Sørlie T, Eisen MB, Van de Rijn M, Jeffrey SS, Rees CA, Pollack JR, Ross DT, Johnsen H, Akslen LA, Fluge O, Pergamenschikov A, Williams C, Zhu SX, Lønning PE, Børresen-Dale AL, Brown PO, Botstein D. Molecular portraits of human breast tumours. *Nature*. 2000;406(6797):747–52. <https://doi.org/10.1038/35021093>.
- Singh D, Febbo PG, Ross K, Jackson DG, Manola J, Ladd C, Tamayo P, Renshaw AA, D'Amico AV, Richie JP, Lander ES, Loda M, Kantoff PW, Golub TR, Sellers WR. Gene expression correlates of clinical prostate cancer behavior. *Cancer Cell*. 2002;1(2):203–9. [https://doi.org/10.1016/S1535-6108\(02\)00030-2](https://doi.org/10.1016/S1535-6108(02)00030-2).

3. van 't Veer LJ, Dai H, Van de Vijver MJ, He YD, Hart AAM, Mao M, Peterse HL, Van der Kooy K, Marton MJ, Witteveen AT, Schreiber GJ, Kerkhoven RM, Roberts C, Linsley PS, Bernards R, Friend SH. Gene expression profiling predicts clinical outcome of breast cancer. *Nature*. 2002;415(6871):530–6. <https://doi.org/10.1038/415530a>.
4. D'Amico AV, Whittington R, Malkowicz SB, Schultz D, Blank K, Broderick GA, Tomaszewski JE, Renshaw AA, Kaplan I, Beard CJ, et al. Biochemical outcome after radical prostatectomy, external beam radiation therapy, or interstitial radiation therapy for clinically localized prostate cancer. *Jama*. 1998;280(11):969–74.
5. Bibikova M, Chudin E, Arsanjani A, Zhou L, Garcia EW, Modder J, Kostelec M, Barker D, Downs T, Fan JB, Wang-Rodriguez J. Expression signatures that correlated with Gleason score and relapse in prostate cancer. *Genomics*. 2007;89(6):666–72. <https://doi.org/10.1016/j.ygeno.2007.02.005>.
6. Penney KL, Sinnott JA, Fall K, Pawitan Y, Hoshida Y, Kraft P, Stark JR, Fiorentino M, Perner S, Finn S, et al. mRNA expression signature of gleason grade predicts lethal prostate cancer. *J Clin Oncol*. 2011;29(17):2391.
7. Sinnott JA, Peisch SF, Tyekucheva S, Gerke T, Lis R, Rider JR, Fiorentino M, Stampfer MJ, Mucci LA, Loda M, et al. Prognostic utility of a new mRNA expression signature of gleason score. *Clin Cancer Res*. 2017;23(1):81–87.
8. Jhun MA, Geybels MS, Wright JL, Kolb S, April C, Bibikova M, Ostrander EA, Fan J-B, Feng Z, Stanford JL. Gene expression signature of gleason score is associated with prostate cancer outcomes in a radical prostatectomy cohort. *Oncotarget*. 2017;8(26):43035.
9. Latil A, Bièche I, Chêne L, Laurendeau I, Berthon P, Cussenot O, Vidaud M. Gene expression profiling in clinically localized prostate cancer: a four-gene expression model predicts clinical behavior. *Clin Cancer Res*. 2003;9(15):5477–85.
10. Long Q, Xu J, Osunkoya AO, Sannigrahi S, Johnson BA, Zhou W, Gillespie T, Park JY, Nam RK, Sugar L, Stanimirovic A, Seth AK, Petros JA, Moreno CS. Global transcriptome analysis of formalin-fixed prostate cancer specimens identifies biomarkers of disease recurrence. *Cancer Res*. 2014;74(12):3228–37. <https://doi.org/10.1158/0008-5472.CAN-13-2699>.
11. Ren S, Wei G-H, Liu D, Wang L, Hou Y, Zhu S, Peng L, Zhang Q, Cheng Y, Su H, et al. Whole-genome and transcriptome sequencing of prostate cancer identify new genetic alterations driving disease progression. *Eur Urol*. 2018;73(3):322–39.
12. Sinha A, Huang V, Livingstone J, Wang J, Fox NS, Kurganovs N, Ignatchenko V, Fritsch K, Donmez N, Heisler LE, et al. The proteogenomic landscape of curable prostate cancer. *Cancer Cell*. 2019;35(3):414–27.
13. Erho N, Crisan A, Vergara IA, Mitra AP, Ghadessi M, Buerki C, Bergstralh EJ, Kollmeyer T, Fink S, Haddad Z, et al. Discovery and validation of a prostate cancer genomic classifier that predicts early metastasis following radical prostatectomy. *PLoS ONE*. 2013;8(6):66855.
14. Karnes RJ, Bergstralh EJ, Davicioni E, Ghadessi M, Buerki C, Mitra AP, Crisan A, Erho N, Vergara IA, Lam LL, Carlson R, Thompson DJS, Haddad Z, Zimmermann B, Sierocinski T, Triche TJ, Kollmeyer T, Ballman KV, Black PC, Klee GG, Jenkins RB. Validation of a genomic classifier that predicts metastasis following radical prostatectomy in an at risk patient population. *J Urol*. 2013;190(6):2047–53. <https://doi.org/10.1016/j.juro.2013.06.017>.
15. Klein EA, Yousefi K, Haddad Z, Choeurung V, Buerki C, Stephenson AJ, Li J, Kattan MW, Magi-Galluzzi C, Davicioni E. A genomic classifier improves prediction of metastatic disease within 5 years after surgery in node-negative high-risk prostate cancer patients managed by radical prostatectomy without adjuvant therapy. *Eur Urol*. 2015;67(4):778–86. <https://doi.org/10.1016/j.eururo.2014.10.036>.
16. Shahabi A, Lewinger JP, Ren J, April C, Sherrod AE, Hacia JG, Daneshmand S, Gill I, Pinski JK, Fan J-B, Stern MC. Novel gene expression signature predictive of clinical recurrence after radical prostatectomy in early stage prostate cancer patients. *Prostate*. 2016;76(14):1239–56. <https://doi.org/10.1002/pros.23211>.
17. Eggener SE, Rumble RB, Armstrong AJ, Morgan TM, Crispino T, Cornford P, Van der Kwast T, Grignon DJ, Rai AJ, Agarwal N, Klein EA, Den RB, Beltran H. Molecular biomarkers in localized prostate cancer: ASCO guideline. *J Clin Oncol*. 2020;38(13):1474–94. <https://doi.org/10.1200/JCO.19.02768>.
18. Morillon A, Gautheret D. Bridging the gap between reference and real transcriptomes. *Genome Biol*. 2019;20(1):1–7.
19. Audoux J, Philippe N, Chikhi R, Salson M, Gallopin M, Gabriel M, Le Coz J, Drouineau E, Commes T, Gautheret D. DE-kupl: exhaustive capture of biological variation in RNA-seq data through k-mer decomposition. *Genome Biol*. 2017;18(1):243. <https://doi.org/10.1186/s13059-017-1372-2>.
20. Pinskaya M, Saci Z, Gallopin M, Gabriel M, Nguyen HTN, Firlej V, Describes M, Rapinat A, Gentien D, De La Taille A, Londoño-Vallejo A, Allory Y, Gautheret D, Morillon A. Reference-free transcriptome exploration reveals novel RNAs for prostate cancer diagnosis. *Life Sci Alliance*. 2019;2(6):1–12. <https://doi.org/10.26508/lsa.201900449>.
21. Thomas A, Barriere S, Broseus L, Brooke J, Lorenzi C, Villemain J.-p., Beurier G, Sabatier R, Reynes C, Mancheron A, Ritchie W. GECKO is a genetic algorithm to classify and explore high throughput sequencing data. *Commun Biol*. 2019;2(1):222. <https://doi.org/10.1038/s42003-019-0456-9>.
22. Michiels S, Koscielny S, Hill C. Prediction of cancer outcome with microarrays: a multiple random validation strategy. *Lancet*. 2005;365(9458):488–92. [https://doi.org/10.1016/S0140-6736\(05\)17866-0](https://doi.org/10.1016/S0140-6736(05)17866-0).
23. Ein-Dor L, Zuk O, Domany E. Thousands of samples are needed to generate a robust gene list for predicting outcome in cancer. *Proc Natl Acad Sci*. 2006;103(15):5923–8. <https://doi.org/10.1073/pnas.0601231103>.
24. Michiels S, Koscielny S, Hill C. Interpretation of microarray data in cancer. *Br J Cancer*. 2007;96(8):1155–8. <https://doi.org/10.1038/sj.bjc.6603673>.
25. Venet D, Dumont JE, Detours V. Most random gene expression signatures are significantly associated with breast cancer outcome. *PLoS Comput Biol*. 2011;7(10):1002240. <https://doi.org/10.1371/journal.pcbi.1002240>.
26. Abeshouse A, Ahn J, Akbani R, Ally A, Amin S, Andry CD, Annala M, Aprikian A, Armenia J, Arora A, et al. The molecular taxonomy of primary prostate cancer. *Cell*. 2015;163(4):1011–25.
27. Fraser M, Sabelnykova VY, Yamaguchi TN, Heisler LE, Livingstone J, Huang V, Shiah Y-J, Yousif F, Lin X, Masella AP, et al. Genomic hallmarks of localized, non-indolent prostate cancer. *Nature*. 2017;541(7637):359–64.
28. Stelloo S, Nevedomskaya E, Kim Y, Schuurman K, Valle-Encinas E, Lobo J, Krijgsman O, Peeper DS, Chang SL, Feng FY-C, et al. Integrative epigenetic taxonomy of primary prostate cancer. *Nat Commun*. 2018;9(1):1–12.
29. Liu J, Lichtenberg T, Hoadley KA, Poisson LM, Lazar AJ, Cherniack AD, Kovatich AJ, Benz CC, Levine DA, Lee AV, et al. An integrated tcga pan-cancer clinical data resource to drive high-quality survival outcome analytics. *Cell*. 2018;173(2):400–16.
30. Ambrose C, McLachlan GJ. Selection bias in gene extraction on the basis of microarray gene-expression data. *Proc Natl Acad Sci*. 2002;99(10):6562–6. <https://doi.org/10.1073/pnas.102102699>.
31. Marçais G, Kingsford C. A fast, lock-free approach for efficient parallel counting of occurrences of k-mers. *Bioinformatics*. 2011;27(6):764–70.
32. Bray NL, Pimentel H, Melsted P, Pachter L. Near-optimal probabilistic rna-seq quantification. *Nat Biotechnol*. 2016;34(5):525–7.
33. Curtin RR, Edel M, Lozhnikov M, Mentekidis Y, Ghaisas S, Zhang S. mpack 3: a fast, flexible machine learning library. *J Open Source Softw*. 2018;3:726. <https://doi.org/10.21105/joss.00726>.
34. Friedman J, Hastie T, Tibshirani R. Regularization paths for generalized linear models via coordinate descent. *J Stat Softw*. 2010;33(1):1–22.
35. Meinshausen N, Bühlmann P. Stability selection. *J R Stat Soc Ser B Stat Methodol*. 2010;72(4):417–73. <https://doi.org/10.1111/j.1467-9868.2010.00740.x>.
36. Kuhn M. Building predictive models in r using the caret package. *J Stat Softw Artic*. 2008;28(5):1–26. <https://doi.org/10.18637/jss.v028.i05>.
37. Menardi G, Torelli N. Training and assessing classification rules with imbalanced data. *Data Min Knowl Disc*. 2014;28(1):92–122.
38. Saito T, Rehmsmeier M. The precision-recall plot is more informative than the ROC plot when evaluating binary classifiers on imbalanced datasets. *PLOS ONE*. 2015;10(3):0118432. <https://doi.org/10.1371/journal.pone.0118432>.
39. Rainer J. EnsDb.Hsapiens.v79: Ensembl based annotation package. R package version 2.99.0. 2017.
40. Liu X, Grogan TR, Hieronymus H, Hashimoto T, Mottahedeh J, Cheng D, Zhang L, Huang K, Stoyanova T, Park JW, et al. Low cd38 identifies progenitor-like inflammation-associated luminal cells that can initiate human prostate cancer and predict poor outcome. *Cell Rep*. 2016;17(10):2596–606.
41. Wang T, Liu Z, Guo S, Wu L, Li M, Yang J, Chen R, Xu H, Cai S, Chen H, et al. The tumor suppressive role of camk2n1 in castration-resistant prostate cancer. *Oncotarget*. 2014;5(11):3611.

42. Liu J, Shen J-X, Wu H-T, Li X-L, Wen X-F, Du C-W, Zhang G-J. Collagen 1a1 (col1a1) promotes metastasis of breast cancer and is a potential therapeutic target. *Discov Med*. 2018;25(139):211–23.
43. Wu X, Wang H, Lian Y, Chen L, Gu L, Wang J, Huang Y, Deng M, Gao Z, Huang Y. Gtse1 promotes cell migration and invasion by regulating emt in hepatocellular carcinoma and is associated with poor prognosis. *Sci Rep*. 2017;7(1):1–12.
44. Chen C-L, Mahalingam D, Osmulski P, Jadhav RR, Wang C-M, Leach RJ, Chang T-C, Weitman SD, Kumar AP, Sun L, et al. Single-cell analysis of circulating tumor cells identifies cumulative expression patterns of emt-related genes in metastatic prostate cancer. *Prostate*. 2013;73(8): 813–26.
45. Robinson JT, Thorvaldsdóttir H, Winckler W, Guttman M, Lander ES, Getz G, Mesirov JP. Integrative genomics viewer. *Nat Biotechnol*. 2011;29(1):24–6. <https://doi.org/10.1038/nbt.1754>.
46. Xie W, Xiao H, Luo J, Zhao L, Jin F, Ma J, Li J, Xiong K, Chen C, Wang G. Identification of low-density lipoprotein receptor class a domain containing 4 (ldlr4) as a prognostic indicator in primary gastrointestinal stromal tumors. *Curr Probl Cancer*. 2020;44(6):100593.
47. Mo S, Zhang L, Dai W, Han L, Wang R, Xiang W, Wang Z, Li Q, Yu J, Yuan J, et al. Antisense lncrna ldlrad4-as1 promotes metastasis by decreasing the expression of ldlrad4 and predicts a poor prognosis in colorectal cancer. *Cell Death Dis*. 2020;11(2):1–16.
48. Chen CD, Welsbie DS, Tran C, Baek SH, Chen R, Vessella R, Rosenfeld MG, Sawyers CL. Molecular determinants of resistance to antiandrogen therapy. *Nat Med*. 2004;10(1):33–39.
49. Chen W-Y, Tsai Y-C, Yeh H-L, Suau F, Jiang K-C, Shao A-N, Huang J, Liu Y-N. Loss of spdef and gain of tgfb1 activity after androgen deprivation therapy promote emt and bone metastasis of prostate cancer. *Sci Signal*. 2017;10(492):6826.
50. Leyten GH, Hessels D, Smit FP, Jannink SA, de Jong H, Melchers WJ. Identification of a candidate gene panel for the early diagnosis of prostate cancer. *Clin Cancer Res*. 2015;21(13):3061–70.
51. Bussemakers MJG, van Bokhoven A, Verhaegh GW, Smit FP, Karthaus HFM, Schalken JA, Debruyne FMJ, Ru N, Isaacs WB. DD3: a new prostate-specific gene, highly overexpressed in prostate cancer. *Cancer Res*. 1999;59(23):5975–9.
52. Koutalellis G, Stravodimos K, Avgeris M, Mavridis K, Scorilas A, Lazaris A, Constantinides C. L-dopa decarboxylase (ddc) gene expression is related to outcome in patients with prostate cancer. *BJU Int*. 2012;110(6b):267–73.
53. Mazzu YZ, Armenia J, Chakraborty G, Yoshikawa Y, Si'Ana AC, Nandakumar S, Gerke TA, Pomerantz MM, Qiu X, Zhao H, et al. A novel mechanism driving poor-prognosis prostate cancer: overexpression of the dna repair gene, ribonucleotide reductase small subunit m2 (rrm2). *Clin Cancer Res*. 2019;25(14):4480–92.
54. Zhong W-D, Liang Y-X, Liang Y-K, Zhuo Y-J, Ye J-H, Zhu X-J, Cai Z-D, Lin Z-Y, Zhu J-G, Wu S-L, et al. Tumor suppressor role and clinical implication of the fifth ewing variant (fev) gene, an ets family gene, in prostate cancer. *In: Prostate Cancer*; 2019. SSRN: <https://ssrn.com/abstract=3372417>.
55. Munkley J, McClurg UL, Livermore KE, Ehrmann I, Knight B, McCullagh P, Mcgrath J, Crundwell M, Harries LW, Leung HY, et al. The cancer-associated cell migration protein tspan1 is under control of androgens and its upregulation increases prostate cancer cell migration. *Sci Rep*. 2017;7(1):1–11.
56. Wan S, Xi M, Zhao H-B, Hua W, Liu Y-L, Zhou Y-L, Zhuo Y-J, Liu Z-Z, Cai Z-D, Wan Y-P, et al. Hmgcs2 functions as a tumor suppressor and has a prognostic impact in prostate cancer. *Pathol Res Pract*. 2019;215(8): 152464.
57. Klein EA, Cooperberg MR, Magi-Galluzzi C, Simko JP, Falzarano SM, Maddala T, Chan JM, Li J, Cowan JE, Tsiatis AC, Cherbavaz DB, Pelham RJ, Tenggara-Hunter I, Baehner FL, Knezevic D, Febbo PG, Shak S, Kattan MW, Lee M, Carroll PR. A 17-gene assay to predict prostate cancer aggressiveness in the context of gleason grade heterogeneity, tumor multifocality, and biopsy undersampling. *Eur Urol*. 2014;66(3):550–60. <https://doi.org/10.1016/j.eururo.2014.05.004>.
58. Witten DM, Tibshirani R. Survival analysis with high-dimensional covariates. *Stat Methods Med Res*. 2010;19(1):29–51. <https://doi.org/10.1177/0962280209105024>.
59. de Ronde JJ, Rigai G, Rottenberg S, Rodenhuis S, Wessels LFA. Identifying subgroup markers in heterogeneous populations. *Nucleic Acids Res*. 2013;41(21):200. <https://doi.org/10.1093/nar/gkt845>.
60. Campos-Laborie FJ, Riusueño A, Ortiz-Estévez M, Rosón-Burgo B, Droste C, Fontanillo C, Loos R, Sánchez-Santos JM, Trotter MW, De Las Rivas J.

DECO: decompose heterogeneous population cohorts for patient stratification and discovery of sample biomarkers using omic data profiling. *Bioinformatics*. 2019;35(19):3651–62. <https://doi.org/10.1093/bioinformatics/btz148>. <https://academic.oup.com/bioinformatics/article-pdf/35/19/3651/30061524/btz148.pdf>.

### Publisher's Note

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

**Ready to submit your research? Choose BMC and benefit from:**

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

**At BMC, research is always in progress.**

Learn more [biomedcentral.com/submissions](https://biomedcentral.com/submissions)

