

# TIF-Seq2 disentangles overlapping isoforms in complex human transcriptomes

Jingwen Wang<sup>1,†</sup>, Bingnan Li<sup>1,†</sup>, Sueli Marques<sup>1</sup>, Lars M. Steinmetz<sup>2,3,4</sup>, Wu Wei<sup>3,5,6</sup> and Vicent Pelechano<sup>1,\*</sup>

<sup>1</sup>SciLifeLab, Department of Microbiology, Tumor and Cell Biology. Karolinska Institutet, Solna, Sweden, <sup>2</sup>European Molecular Biology Laboratory (EMBL), Genome Biology Unit, Heidelberg, Germany, <sup>3</sup>Stanford Genome Technology Center, Stanford University, Palo Alto, CA, USA, <sup>4</sup>Department of Genetics, School of Medicine, Stanford University, Stanford, CA, USA, <sup>5</sup>CAS Key Laboratory of Computational Biology, CAS-MPG Partner Institute for Computational Biology, Shanghai Institute of Nutrition and Health, Shanghai Institutes for Biological Sciences, University of Chinese Academy of Sciences, Chinese Academy of Sciences, Shanghai 200031, China and <sup>6</sup>Center for Biomedical Informatics, Shanghai Engineering Research Center for Big Data in Pediatric Precision Medicine, Shanghai Children's Hospital, Shanghai Jiao Tong University, Shanghai 200040, China

Received May 12, 2020; Revised July 17, 2020; Editorial Decision August 06, 2020; Accepted August 07, 2020

## ABSTRACT

Eukaryotic transcriptomes are complex, involving thousands of overlapping transcripts. The interleaved nature of the transcriptomes limits our ability to identify regulatory regions, and in some cases can lead to misinterpretation of gene expression. To improve the understanding of the overlapping transcriptomes, we have developed an optimized method, TIF-Seq2, able to sequence simultaneously the 5' and 3' ends of individual RNA molecules at single-nucleotide resolution. We investigated the transcriptome of a well characterized human cell line (K562) and identified thousands of unannotated transcript isoforms. By focusing on transcripts which are challenging to be investigated with RNA-Seq, we accurately defined boundaries of lowly expressed unannotated and read-through transcripts putatively encoding fusion genes. We validated our results by targeted long-read sequencing and standard RNA-Seq for chronic myeloid leukaemia patient samples. Taking the advantage of TIF-Seq2, we explored transcription regulation among overlapping units and investigated their crosstalk. We show that most overlapping upstream transcripts use poly(A) sites within the first 2 kb of the downstream transcription units. Our work shows that, by paring the 5' and 3' end of each RNA, TIF-Seq2 can improve the annotation of complex genomes, facilitate accurate assignment of promoters to genes and easily identify transcriptionally fused genes.

## INTRODUCTION

Eukaryotic transcriptomes are complex, involving thousands of coding and non-coding RNA isoforms differing in transcription start sites (TSSs), poly(A) sites (PASs) and splicing. Overlapping isoforms can have divergent functional consequences: changing the encoded protein (1,2) or affecting mRNA post-transcriptional life (e.g. translation, localization and stability) (3). However, the interleaved nature of the transcriptomes convolutes its study and limits the accurate identification and quantification of alternative isoforms (4,5). This can lead to incomplete or inaccurate annotations, which cause misinterpretation of gene expression data (6) and limit our ability to link regulatory regions with genes (7) (and thus to genetically manipulate them). The correct identification of transcription boundaries of overlapping isoforms is particularly challenging (8), even if we know that alternative TSS and PAS drive most isoform's variations across human tissues (9). Standard RNA-Seq can identify transcribed regions and splicing events, however, it cannot distinguish RNA fragments originating from alternative overlapping features. The overlapping nature of the transcriptomes also limits our ability to dissect the molecular mechanism underlying the regulatory crosstalk across adjacent transcription units (10,11). For example, RNA-Seq cannot distinguish if a specific mRNA fragment in the 5' region of a gene will reach the canonical poly(A) site or originates from an overlapping transcription unit that terminates prematurely. Single-end approaches such as CAGE or poly(A) site sequencing have been key to define the boundaries of the transcriptomes (12,13). However, those approaches cannot study the combination between TSS and PAS. Long-read sequencing technologies promise to reveal

\*To whom correspondence should be addressed. Tel: +46 72 8564904; Email: vicente.pelechano.garcia@ki.se

†The authors wish it to be known that, in their opinion, the first two authors should be regarded as Joint First Authors.

transcription complexity on the genome-wide scale (14). Nevertheless, the high cost combined with low throughput and limited resolution in the 5' and 3' transcript regions are still major limitations (15).

To bridge the gap between short-read and long-read technologies, and to improve our ability to study the regulatory crosstalk between overlapping transcription units on the same strand, we have developed an optimized Transcript Isoform Sequencing (TIF-Seq2) that is especially well-suited for the interrogation of complex transcriptomes. TIF-Seq2 allows to sequence simultaneously the 5' and 3' ends of individual RNA molecules at single-nucleotide resolution. To demonstrate its utility, we dissected the overlapping transcriptome of a chronic myeloid leukaemia (CML) cell line (i.e. K562) in response to imatinib treatment. We identified thousands of known and unannotated transcript isoforms, accurately defined the boundaries of lowly expressed intergenic transcripts and validated them using alternative short-read and targeted long-read sequencing approaches. We focused on overlapping transcription units that are particularly challenging to investigate with RNA-Seq, and showed the common existence of short overlapping upstream transcripts that may lead to misinterpretation of RNA-Seq and CAGE gene expression data. We also showed the existence of more complex overlapping and read-through transcripts. Finally, we used the obtained information to improve the detection and analysis of complex overlapping transcripts in clinical RNA-Seq datasets and showed evidences of transcriptional events involving gene-promoter rewiring and potentially leading to the generation of transcriptionally fused proteins.

## MATERIALS AND METHODS

### Cell culture

The human erythroleukemia cell line K562 was obtained from ATCC(ATCC®CCL-243™). Cells were cultured in RPMI 1640 medium supplemented with 10% FBS, 2 mM L-glutamine, 1% pen/strep (Life Technologies) at 37°C with 5% CO<sub>2</sub> in a humidified atmosphere. Cells (3 × 10<sup>5</sup> cells/ml) were exposed to 0.2–5 μM (0.2, 0.5, 1, 1.7 and 5 μM) imatinib for 8, 24 and 48 h. Aliquots were taken at each time point for the assessment of cell viability via Trypan blue staining by EVE™ automated cell counter. Two biological replicates of K562 cells treated with 1 μM imatinib for 24 h and corresponding DMSO control were used for TIF-Seq2 library preparation.

### TIF-Seq2 library preparation

In brief, capped and polyadenylated RNA was used as a template to generate full-length cDNA. We circularized the cDNA, removed non circularized molecules and fragmented circularized molecules using sonication. Streptavidin magnetic beads were used to purify the fragments spanning the 5' and 3' ends of cDNA and then were used for Illumina library preparation (Supplementary Figure S1). In detail, 2.5 μg total RNA was treated with Turbo DNase (0.12 U/μl) (Fisher Scientific) for 20 min at 37°C to prevent genomic DNA contamination. After inactivation of DNase, input RNA was dephosphorylated

by incubating with Calf Intestinal Alkaline Phosphatase (0.3U/μl) (CIP) (NEB) at 37°C for 30 min. Two rounds of phenol-chloroform extraction followed by ethanol precipitation were performed to remove CIP. After dephosphorylation, input RNA was decapped by incubating with Cap-Clip Acid Pyrophosphatase (0.125 U/μl) (CellScript) at 37°C for 60min. After phenol–chloroform purification and ethanol precipitation, RNA was ligated overnight at 16°C at 5' with DNA/RNA chimeric oligonucleotide adaptor (TCAGACGTGTGCTCTTCCGATCTrNrNrWrNrNrWrNrN, TIF2-RNA in Supplementary Table S1) using T4 RNA ligase (NEB) in the presence of 10% dimethylsulphoxide (DMSO), RiboLock RNase Inhibitor (Thermo Fisher Scientific EO0382) and 1mM ATP. The chimeric TIF2-RNA adaptor introduced a common anchor sequence for forward primer of subsequent PCR amplification and an 8-mer unique molecular identifier (UMI). Ligated RNA was purified with 1.8:1 volumetric ratio (1.8×) RNA clean XP beads according to the manufacturer's instruction and then used as the template for reverse transcription (RT). Ligated RNAs were reverse-transcribed using barcoded oligo-dT primers (i.e., TAGTTCAGTCTTCAGTACCTCGTGC GGCCGCX XXXXXACACTCTTTCCTACACGACGCTCTTC CGATCTTTTTTTTTTTTTTTTTTVN; where X refers to the specific barcode, TIF2-RT in Supplementary Table S1) which introduced Illumina sequencing primer 1, a 3' index, a NotI endonuclease digestion site and a 3' common sequence for the subsequent PCR reaction. Specifically, RNAs were mixed with the corresponding TIF2-RT oligo and dNTPs, denatured at 65°C for 5 min and then put on ice. The sample was mixed with 5× First-strand buffer, Trehalose (1.57M) and RiboLock RNase inhibitor and incubated first at 42°C for 2 min. Finally, 2 μl of SuperScript™ III reverse transcriptase (Thermo Fisher Scientific) was added to each reaction and incubated at 42°C for 50 min, at 50°C for 30 min, at 55°C for 30 min and inactivated at 70°C for 15 min. Used RNA template was removed by incubating with 0.5 ul RNase H (5 U/μl) and 0.5 μl RNase cocktail (Ambion) at 37°C for 30 min. First-strand cDNA was purified with 2X Ampure XP beads according to manufacturer instruction. To avoid saturating the reaction, only half of the obtained cDNA was used for the following PCR, and the rest was stored as a backup. Used cDNA template was further split into two PCR reactions with Terra PCR Direct Polymerase (Takara) with the following program 98°C for 2 min, then 16 cycles of 98°C for 20 s, 60°C for 30 s, 68°C for 5 min (+10 s/cycle) and finally 72°C for 5 min. The PCR above used as primers TIF2-Rv: TAGTTCAGTCTTCAGTACCTCGT and TIF2-Fw: TATAGCGCCGCXXXXXXGTGAC[BtdnT]GGAG TTCAGACGTGTGCTCTTCCGATC (where X refers to different barcodes). PCR products were purified with 1X Ampure XP beads according to the manufacturer's instruction and then quantified with Qubit dsDNA HS assay. PCR products from different samples were then pooled together with equal mass. By pooling full-length cDNA containing sample-specific barcodes, we were able to detect and estimate the percentage of intermolecular circularization events (i.e. chimeras connecting barcodes originating from different samples). This information

was used to select the optimal concentration favouring intramolecular ligation (see below). Pooled PCR products were subjected to 1U/ $\mu$ l NotI HF (New England BioLabs) endonuclease digestion at 37°C for 1 h. After inactivation of NotI at 65°C for 20 min, samples were purified with 1.8 $\times$  Ampure XP beads according to the manufacturer's instruction.

To favour intramolecular ligation, PCR products with sticky ends were highly diluted to a final concentration less than 1 ng/ $\mu$ l and ligated with a high concentration (66.68 U/ $\mu$ l) of T4 DNA ligase (New England BioLabs) at 16°C for at least 16 h. To remove the unligated linear PCR products, we added 0.5  $\mu$ l plasmid-safe for every 100  $\mu$ l total volume and incubated at 37°C for one hour in the presence of 1  $\mu$ l ATP 100 mM. After inactivation at 70°C for 30 min, the self-circularized cDNA was purified with phenol-chloroform and ethanol precipitation. Circularized cDNA was fragmented by sonication (Covaris ME220; Covaris, Inc) (duration 240 s, peak power 30, duty factor 10, 200 cycles/burst, average power 3 W). The fragments were then purified with 1X Ampure XP beads according to manufacturer instruction. Biotin labelled fragments which contained the 5' and 3' connecting region was enriched using Dynabeads M280 streptavidin (Invitrogen) and incubated at room temperature for 30 min. Captured fragments were subjected to end repair with End Repair Enzyme Mix (New England Biolabs) and dA tailing with Klenow Fragment (3' - 5' exo-) 5U/ $\mu$ l (New England Biolab). To add the required Illumina grafting sequences, each sample (20  $\mu$ l of resuspended beads) was incubated at room temperature for one hour with a mixture of 10  $\mu$ l 5 $\times$  Quick ligation buffer (New England BioLabs), 16  $\mu$ l nuclease-free water (Ambion), 3  $\mu$ l T4 DNA ligase(2000U/ $\mu$ l) and 1  $\mu$ l 1 $\mu$ M duplex adaptors. Duplex adaptors were generated by annealing TIF2-forkFw (5'-AATGATACGGCGACCACC GAGATCTACACACACCTGCCGGTCACC\*T-3') and TIF2-forkRv (5'-phos-GGTGACCGGCAGGTGTATCT CGTATGCCGCTTCTGCTTG-3') at 15  $\mu$ M and diluted to 1  $\mu$ M working solution freshly when used each time. After cleaning the beads, the beads resuspended in 20  $\mu$ l EB buffer were used for PCR amplification using 25  $\mu$ l Phusion High-fidelity MasterMix (2 $\times$ ) (New England BioLabs), 4  $\mu$ l nuclease-free water and 0.5  $\mu$ l PCRgraftP5 primer (5 uM, 5'-AATGATACGGCGACCACCAGAGATCTACAC-3') and 0.5  $\mu$ l PCRgraftP7 primer (5 uM, 5'- CAAGCAGAAGACGGCATAACGAGAT-3'). Samples were subjected to PCR amplification with the following program: 98°C for 30 s, then 18 cycles of [98°C for 20 s, 65°C for 30 s, 72°C for 30s ] and finally 72°C for 5 min. DNA library then went through two-step beads-based size selection (first step 0.35 $\times$ , the second step take the supernatant from first step and add extra 0.45 $\times$ ) with expected size distribution in the range of 300–1000 bp. Purified DNA library was adjusted to 4 nM, then denatured, diluted according to Illumina's instruction and sequenced on Nextseq 500 platform.

We used four custom sequencing oligos as follows: SeqR1+15T (5'-ACACTCTTCCCTACACGACGCTC TTCCGATCTTTTTTTTTTTTTTTT-3'), SeqINDX1 (5'-GATCGGAAGAGCACACGTCTGAACTCCAGT CAC-3'), SeqINDX2 (5'-GATCGGAAGAGCGTCG

TGTAGGGAAAGAGTGT-3') and SeqR2 (5'-GTGA CTGGAGTTCAGACGTGTGCTCTTCCGATCT-3'). See Supplementary Figure S2 for details regarding the sequencing oligonucleotide annealing. Two mixtures of sequencing oligos were prepared as follows: mixture 1 contained SeqR1+15T (0.3  $\mu$ M) and SeqR2 (0.3  $\mu$ M), mixture 2 composed of SeqINDX1 primer (0.3  $\mu$ M) and SeqINDX2 primer (0.3  $\mu$ M). Mixture 1 was loaded into both positions (#7 and #8 on reagent cartridge in a NextSeq 500 instrument) for custom read1 primer and for custom read2 primer. Mixture 2 was loaded into the position for custom index primer (#9 on reagent cartridge). Sequencing was carried out in an Illumina NextSeq 500 instrument with stand-alone configuration and custom sequencing oligos. Paired-end sequencing read lengths were set as read1 76 bp, read2 76 bp, index1 6 bp and index2 6 bp.

### 3'T-fill library preparation

We performed 3'T-fill as previously described (16). In brief, 10  $\mu$ g DNA-free total RNA (16  $\mu$ l) was fragmented by adding 4  $\mu$ l fragmentation buffer (5 $\times$ ) (200 mM Tris-acetate, pH 8.1, 500 mM potassium acetate, 150 mM magnesium acetate) and incubated at 80°C for 5 min. Fragmented RNA was purified with 1.5 $\times$  Ampure XP beads according to the manufacturer's instruction and used as the template for reverse transcription (RT). The sample was mixed with biotinylated dT primer (P5\_dT16VN 5'-[BtN]AATGATACGGCGACCACCAGATCTACACT CTTTCCCTACACGACGCTCTTCCGATCTTTTTTT TTTTTTTTTVN-3') (where V refers to A, C or G) and denatured at 65°C for 5 min. Denatured RNA was then mixed with 4  $\mu$ l 5 $\times$  first strand buffer, 2  $\mu$ l DTT (0.1 M) and 4  $\mu$ l freshly prepared actinomycin D(0.1  $\mu$ g/ $\mu$ l). After incubation at 42°C for 2 min, 0.5  $\mu$ l SuperScript™ II (Invitrogen) was added, incubated at 42°C for 50 min and finally 72°C for 15 min. cDNA was purified with 1.5 $\times$  Ampure XP beads according to the manufacturer's instruction. For second-strand synthesis, cDNA was mixed with 0.5  $\mu$ l RNase H (5 U/ $\mu$ l) (New England BioLabs) and 2  $\mu$ l DNA Polymerase I (10U/ $\mu$ l) (Thermo Scientific) and incubated at 16°C for 2.5 h. cDNA was purified with 0.9 $\times$  Ampure XP beads according to the manufacturer's instruction. 20  $\mu$ l biotin labelled cDNA was mixed with 20  $\mu$ l Dynabeads M280 streptavidin beads (Invitrogen) at room temperature for 15 min. Following end repair and dA tailing, each sample (8  $\mu$ l) was mixed with 12.5  $\mu$ l 2 $\times$  Quick Ligation buffer (New England BioLabs), 2.5  $\mu$ l T4 DNA ligase (2000 U/ $\mu$ l, New England BioLabs) 2  $\mu$ l of annealed duplexed adaptors at 2.5  $\mu$ M(5'-GTGACTGGAGTTCAGACGTGTGCTCTTCCGATC\*T-3', 5'-Phos-GATCGGAA GAGCACACGTCTGAACTCCAGTCAC[AmC7]-3') and incubated for at 20°C for 20 min. Adaptor ligated libraries were washed and amplified by PCR by adding 25  $\mu$ l 2X Fusion High Fidelity Master Mix, 0.5  $\mu$ l PE1 primer (10  $\mu$ M, 5'-AATGATACGGCGACCACCAGATCTA CACTCTTCCCTACACGACGCTCTTCCGATC\*T-3')and 0.5 $\mu$ l PE2\_MPX primers (10  $\mu$ M, 5'-CAAGCAGAAGACGGCATAACGAGATXXXXXXGT GACTGGAGTTCAGACGTGTGCTCTTCCGATC\*T-



3'; where X refers to the specific barcode). Samples were subjected to the following program: 98°C for 30 s, 18 cycles of 98°C for 10 s, 65°C for 10 s, 72°C for 10 s and finally 72°C for 5 min. Library was purified with 1.8× Ampure XP beads according to the manufacturer's instruction.

### PacBio long reads sequencing

To validate the existence and structure of unannotated transcript isoforms, we used the TIF-Seq2 derived information to perform targeted amplification followed by Pacific Biosciences (PacBio) sequencing. To decrease the costs associated to library preparation, we performed a pooled reverse transcription with a mix of gene specific primers and then split the obtained cDNAs for individual PCR reactions. We used two biological replicates of K562 exposed or not to imatinib (four samples, as in the TIF-Seq2) as input. Each sample was individually labelled during the PCR reaction (see below). We pooled equal amount of all isoform-specific RT primers (Supplementary Table S1) and used total RNA as a template for reverse transcription (RT). Isoform-specific RT primers (PB\_RT) were designed with a universal primer sequence (5'-GTGACTGG AGTTCAGACGTGT), plus eight random nucleotides as unique molecular identifier and a gene specific sequence. RNA was mixed with primer and dNTPs, denatured at 65°C for 5 min and transferred to ice. For the reverse transcription, we added first strand buffer, Trehalose (1.57 M), Ribolock RNase Inhibitor and the mix was incubated at 42°C for 2 min, and then 2 µl of SuperScript II reverse transcriptase (ThermoFisher Scientific) was added to each reaction. We incubated each reverse transcription reaction at 42°C for 50 min, 50°C for 30 min, 55°C for 30 min and inactivated at 70°C for 15 min. Obtained cDNA was distributed in 96-well plates containing in each well an isoform-specific forward primer (PB\_FW, composed of a common sequence 5'-ACACTCTTCCCTACACGAC and gene specific sequence) and a common reverse primer containing 8 mer sample identifier (Supplementary Figure 12). cDNA was PCR amplified using the Phusion High Fidelity Master Mix with the following program: denaturation at 98°C for 2 min, then 98°C 20 s, 60°C 30 s, 68°C 5 min (+10 s/cycle) for 35 cycles and final extension was performed at 72°C for 5 min. Amplified products were analysed by gel electrophoresis, purified by 1.8× Ampure XP beads, and pooled at similar concentration for sequencing. Pooled PCR products were used for library generation using SMRTbell™ Template Prep Kit 1.0-SPv3 and sequenced on PacBio Sequel system.

### TIF-Seq and 3'T-fill sequencing processing and alignment

We employed bcl2fastq (v2.20.0) for converting raw images to sequence information (FASTQ files) and demultiplex, allowing two mismatches in index 1 and one mismatch in index 2. For TIF-Seq2 data, we collapsed all 5'-end or 3'-end sequence reads in each sample according to the indexes. Cutadapt (17) (v1.16) was utilized to trim TIF-seq2 sequencing primer (-a AGGTGACCGGCAGGTGT) and Illumina TruSeq adapter (-a AGATCGGAAG). After extracting 8 bp of unique molecular identifiers (UMIs) with

UMI-tools (18) (v0.5.4) from the 5' ends and removing extra A stretches in the 3' ends caused by poly(A) slippage during PCR amplification, we kept reads over 20 bp for alignment. We used STAR (19) (v2.5.3a) for aligning 5'-end reads and 3'-end reads separately to the human reference genome hg38, supplying Gencode v27 transcripts as splicing junction annotation. Alignment setting was adjusted as below, `-alignIntronMax 200000 -alignEndsType Extend5pOfRead1 -alignSJoverhangMin 10`. We then linked paired-end reads and kept the uniquely mapped pairs that are on the same chromosome (Supplementary Figure S4). Furthermore, a customised script adapted from UMI-tools was employed to remove PCR duplicates from the leftover reads, allowing 1 bp mismatch in the UMIs and 1 bp shifting in the transcription start sites (<https://github.com/jingwen/TIFseq2/blob/master/dedup.py>).

For 3'T-fill sequencing data, we trimmed Illumina TruSeq adapter (-a AGATCGGAAG) and extra A stretches in the 3' ends with Cutadapt v1.16. Reads over 20 bp are aligned to hg38 by using STAR v2.5.3a in paired-end mode, supplying Gencode v27 transcripts as splicing junction annotation with adjusted setting, `-alignSJDBoverhangMin 1 -alignIntronMax 1000000 -alignMatesGapMax 1000000 -alignEndsType Extend5pOfReads12`. Only uniquely mapped reads were kept for downstream analysis.

In order to evaluate the improvement of TIF-Seq2, we compared our data from the current study to a previous study on human HeLa cells using TIF-Seq1 (20) (GEO: GSE75183). We aligned TIF-Seq1 reads to the human reference genome hg38 using STAR with the same parameter setting as how we analysed TIF-Seq2 data.

### Transcription boundary determination

A customised python script was employed to extract both boundaries of TIF-Seq2 read pairs and 3' tags of 3'T-fill sequencing data, collapse the boundary tags and calculate the coverage (<https://github.com/jingwen/TIFseq2/blob/master/boundary.py>). In order to filter out false positive poly(A) sites caused by internal priming, we performed motif enrichment analysis near putative poly(A) sites according to the composition of adenine (A) in their downstream sequence. Poly(A) associated hexamers (A[AT]TAAA) were discovered at 15–30 nt upstream of poly(A) sites with less than 7 As, while no obvious motif was detected near the 3'-end tags with at least 7 As in the downstream (Supplementary Figure S18 B–D). Therefore, we regarded the 3'-end tags with at least 7 As in the downstream 10 nt sequences as internal priming cases, and excluded them from downstream analysis. Then we employed CAGER (21) to define the cluster of transcripts 5'- or 3'-end tags of TIF-Seq2 respectively (Supplementary Figure S4B). Transcription boundary tags were normalized to match a power-law distribution (22). In brief, the reverse-cumulative distribution of the number of boundary tags with at least a give number of tags were fitted to a common reference power-law. Low-coverage tags supported by less than 1 normalized counts in more than one sample were excluded before clustering. The boundary tags within 10 bp window were spatially clustered together. Clusters with only one boundary tag are kept if the normalized counts are above 1.

The tag clusters were further formed into non-overlapping consensus clusters across all samples if they are within 10 bp apart. The same strategy was applied for identifying consensus PAS clusters from 3'T-fill sequencing data.

### TIF definition, annotation and quantification

We then linked the 5'-end TSS and 3'-end PAS clusters according to the supporting read pairs from TIF-Seq2 (Supplementary Figure S4B). We filtered out pairs with extremely long (>2 Mb) and extremely short (<300 bp) mate-pair distance. To keep a conservative estimate of unannotated transcript isoforms identified, we excluded pairs mapping to different chromosomes. Transcript isoform boundaries (TIFs) were defined as connection between TSS clusters and PAS clusters supported by at least four read pairs connecting them across all samples (unique molecular events). The TIFs were further assigned to Gencode v28 annotation features based on their relative distance to the annotated transcripts (Supplementary Figure S4C). TSS distances (d1) and PAS distances (d2) were calculated between a TIF and its overlapping annotated transcripts. A TIF was assigned to the transcript with the least sum of d1 and d2 among all overlapping transcripts, further assigned to the gene that harbours the transcript. According to the relative position to their assigned transcripts, the TIFs were classified as (i) annotated transcripts, if both TSS and PAS are within 200 bp away of annotated transcripts boundaries; (ii) transcripts with new TSS; (iii) transcripts with new PAS; (iv) transcripts with new boundaries, if both TSS and PAS not annotated and (v) intergenic TIFs (Supplementary Figure S9). We measured TIF-Seq2 expression in K562 cells as count of read pairs that link TIF boundaries. We employed DESeq2 (23) for normalization and differential expression analysis (before and after drug treatment) with default setting.

### Long-read sequencing data analysis

We trimmed the PCR primers from both ends of the PacBio highly accurate consensus sequences using Cutadapt (v1.16) and extracted UMIs with UMI-tools (v1.0.0), keeping the reads with at least 200 bp in length. Then the reads were aligned to human reference genome hg38 using minimap2 (24) (v2.16) with the following setting (-ax splice -uf -C5 -O6,24 -B4). We used UMI-tools for removing the PCR duplicates with adjusted setting as -method cluster -spliced-is-unique. We further employed BEDTools (25) (v2.27.1) bamToBed function to convert the alignment reads into BED format and a customised script to convert BED format to GFF format.

### Independent RNA-Seq validation

RNA-Seq data of K562 cells before and after imatinib treatment ( $n = 4$ ) from study Gallipoli *et al.* (26) were downloaded from GEO depositories (GSE105161). RNA-Seq data of 21 paired CML patient (27) before and after imatinib treatment were downloaded from EGA archive (EGAD00001004179). We employed STAR (v2.5.3a) to align paired-end reads to human reference

genome hg38, adding long-read sequencing validated transcripts into splicing junction annotation, with adjusted setting (-alignSJDBoverhangMin 1 -alignIntronMax 200000 -alignEndsType Extend5pOfReads12). Polyribosome profiling data of K562 cells (28) ( $n = 3$ ) were downloaded from GEO (GSE93210). We used HISAT2 (29) (v2.1.0) for aligning the reads to human reference genome hg38 and long-read sequencing validated transcripts as splicing junction annotation, meanwhile adjusting maximum intron length to 200kb.

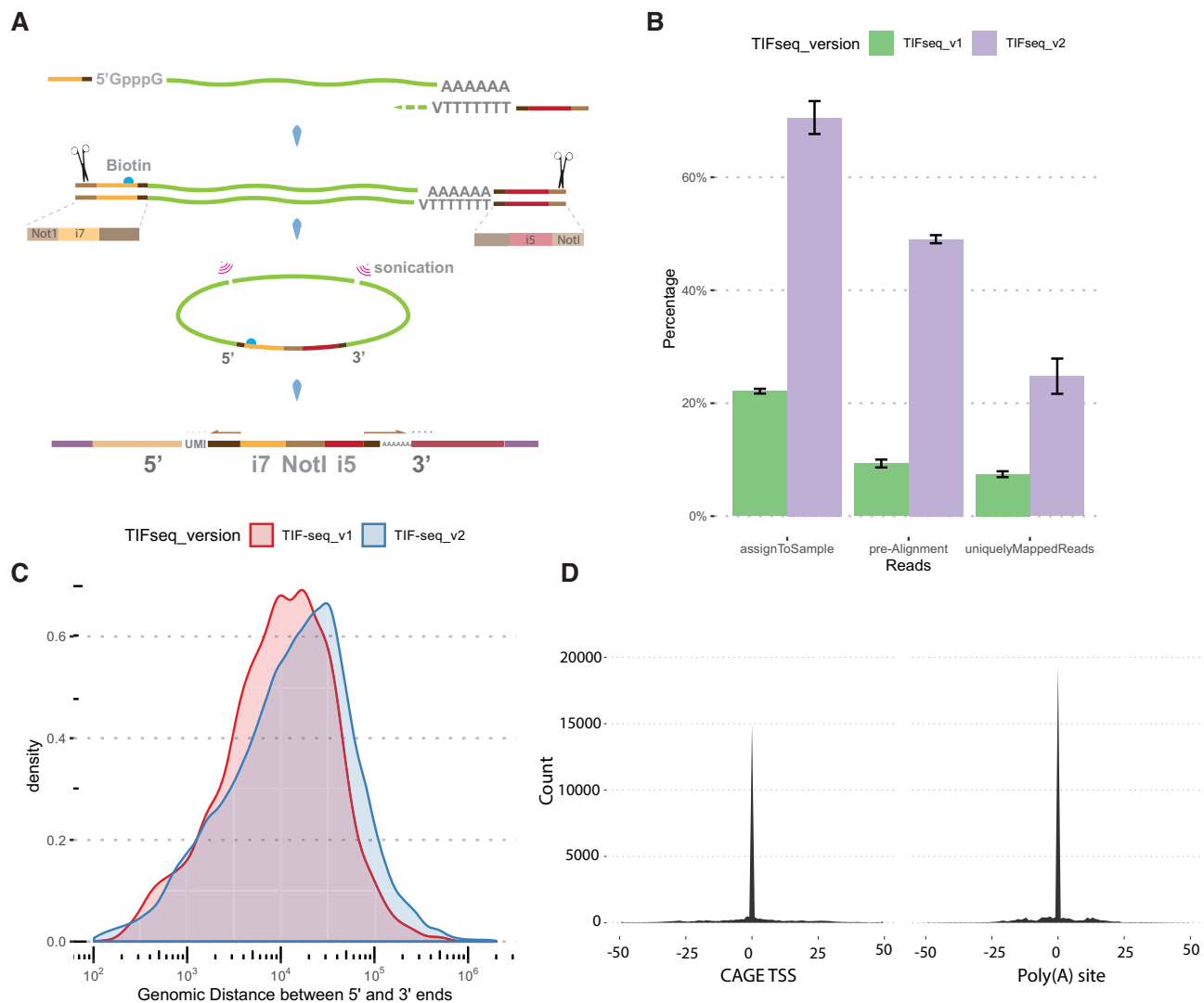
Gene expression in K562 cells (26) and CML patients (27) from standard RNA-Seq was quantified using featureCounts (30) according to Gencode v28 transcript annotation, TIF-Seq2 transcription boundaries and long-read sequencing validated transcript model. We employed DESeq2 (23) for normalization and differential expression analysis (before and after drug treatment) with default setting for K562 cells. For CML patient RNA-Seq data, we set up multiple factors (paired samples, drug treatment and phenotype) in DESeq2 for differential expression test.

## RESULTS

### TIF-Seq2 delineates isoforms in complex human transcriptome

We and others have previously developed approaches able to link the 5' and 3' regions of individual transcripts (31–33). Our work in *S. cerevisiae* demonstrated the existence of a complex overlapping transcriptome even in a simple eukaryote with limited splicing (31). However, the applicability of our original protocol (34) to the study of complex transcriptomes was limited due to the variability in length of the mappable 5' and 3' tags and its modest throughput. To increase the length of the boundary tags, we designed a new sequencing strategy decoupling the region required for bridge amplification from the Illumina sequencing primers (Figure 1A; Supplementary Figures S1 and S2). This decreases the need to perform stringent library size-selection, allows sequencing from the exact 5' and 3' ends of each RNA molecule and generates longer reads suitable for the study of complex genomes (see Methods). In addition, we performed extensive enzymatic optimizations to maximize the length and complexity of the full-length cDNA libraries, as well as introduced early sample pooling, unique molecular identifiers (UMI) and barcodes to control the formation of chimeras (Supplementary Figure S1).

To demonstrate its utility, we investigated the transcriptome of a well-characterized CML cell line (K562) in response to the tyrosine kinase inhibitor imatinib (Supplementary Figure S3). After quality control and PCR deduplication (see Materials and Methods and Supplementary Figure S4A), we obtained over 14 million pair tags uniquely mapped to the human genome at single-nucleotide resolution (Supplementary Tables S2 and S3). Compared with TIF-Seq1 (20), all these modifications improved the number of informative reads, genomic distances of transcription boundaries (Figure 1B and C), and thus allowed the application of TIF-Seq2 to complex genomes. We clustered adjacent transcription start sites (TSSs) and polyadenylation sites (PASs) and obtained 32,631 TSS and 31,187 PAS clusters (see Materials and Methods). Identified clusters are



**Figure 1.** Genome-wide measurement of transcript isoforms with TIF-Seq2. (A) Capped and polyadenylated RNA used as template to generate full-length cDNA, which then circularized and fragmented using sonication. Streptavidin magnetic beads were used to purify the fragments spanning the 5' and 3' end of cDNA and then were used for Illumina library preparation. The arrows indicate the direction of sequencing reads extension (more details in Supplementary Figure S1). (B) Informative reads fetched from TIF-Seq1 ( $n = 3$ ) and TIF-Seq2 ( $n = 4$ ). TIF-Seq2 can fetch more useful reads that are assigned to samples, reads passed quality control for alignment and the uniquely mapped reads. (C) Genomic distance between 5' and 3' ends captured by TIF-Seq1 and TIF-Seq2. The enzymatic optimization of TIF-Seq2 can improve the lengths of RNA molecules (average distance: 20 kb in TIF-Seq1 and 35 kb in TIF-Seq2). (D) Transcript boundaries agree with the transcription start sites (TSSs) defined by CAGE (12) and the poly(A) sites defined by 3' sequencing (13).

narrow, and 80% of them are smaller than 15 nt (TSS) and 10 nt (PAS) (Supplementary Figure S5A and B). In general, the widths of TSS or PAS clusters tend to be associated with gene expression (Spearman correlation coefficients: 0.86 (TSS clusters) and 0.78 (PAS clusters), Supplementary Figure S5C and D). However, we did not observe any positive correlation between the widths of matched TSS and PAS clusters (Supplementary Figure S5E). To validate the accuracy of TIF-seq2, we performed 3'T-fill (16) with the same samples and compared transcript boundaries with independent published datasets for TSS (CAGE (12)) and PAS measurements (13). In both cases, those analyses demonstrate our accurate detection of transcript boundaries (Figure 1D and Supplementary Figure S6). Among

high-confidence (normalised counts > 10) TIF boundaries, 85% of TSS and 92% of PAS clusters are located within 100 bp next to the TSSs and PASs from the published datasets (12,13).

The main advantage of our approach is that it allows to link the identified TSS and PAS clusters (Supplementary Figure S4). Doing so, we identified 49,847 unique combinations of TIF-Seq2-linked TSS-PAS clusters (referred here as **TIFs**, *boundary* transcript isoforms) supported by at least four independent molecular events across all samples. TIF-Seq2-derived transcript boundaries are in good agreement with the curated Gencode v28 annotation (Supplementary Figure S7). At our current sequencing coverage, TIFs overlap 9,006 annotated genes, 80% of which are covered by



more than one TIF (Supplementary Figure S8A). We observed a mild correlation between gene expression and the number of identified TIFs in each gene (Spearman correlation coefficient: 0.43, Supplementary Figure S8B), for highly expressed genes were more likely to be captured while we might identify only the prominent isoforms in lowly expressed genes. In those genes with multiple TIFs, there is usually a dominant TIF with relatively higher expression among all other TIFs harboured in the same gene (Supplementary Figure S8C). In spite of the good agreement (Supplementary Figure S7), 60% of the TIFs support unannotated isoforms with alternative TSSs, PASs, or both (Supplementary Figure S9). Current genome annotation is based mainly on classical cloning strategies, RNA-Seq and CAGE (5,12), which limit its ability to distinguish between overlapping isoforms harboured in the same gene loci. As TIF-Seq2 is able to reveal complex overlapping transcript structures, we decided to focus on those that are often missed by traditional approaches.

### TIF-Seq2 improves transcript boundary delineation in clinical RNA-Seq data

Lowly expressed transcripts are in general challenging to detect. Long-read approaches lack the throughput, and even when combined with capture-based enrichment, they require a prior definition of the regions of interest (35). On the other hand, short-read RNA-Seq approaches with much higher throughput distribute their sequencing power along the whole transcribed region, and need to be linked with independent TSS or PAS dataset to accurately infer their putative boundaries (8). On the contrary, TIF-Seq2 focuses its sequencing power on the transcript boundaries (TSS and PAS), allowing to link them confidently even with relatively low coverage. This simplifies the annotation of lowly expressed transcripts and makes it easy to distinguish them from background noise. Using this approach, we identified 1,034 TIFs in 426 unannotated non-overlapping transcribed regions, defined as unannotated genes (Supplementary Table S4). Although the main application of TIF-Seq2 is the dissection of the overlapping transcription organization, it also allows the quantification of differential expression of the same isoforms across samples. To test that, we treated the K562 cells with 1  $\mu$ M imatinib for 24h. Doing so, we identified 879 TIFs that are up/down-regulated (Supplementary Figure S10 and Supplementary Table S5). And among those, 60 TIFs in 36 unannotated non-overlapping transcripts were differentially expressed (in agreement with our estimation by 3'T-fill, Supplementary Figure S11 and Supplementary Table S6). We found evidence for their transcription and expression pattern using independent RNA-Seq datasets (26) (Supplementary Table S4). About 24% of unannotated transcripts are significantly regulated (adjusted  $P$ -value < 0.05) after imatinib treatment. In order to validate the identified isoforms and determine their internal structure, we used the TIF-Seq2-derived boundaries of 28 differentially expressed unannotated genes for designing amplicon-based enrichment of full-length isoforms followed by long-read sequencing (Supplementary Figure S12). We confirmed the existence of 25 candidates and predicted the existence of short (29-179 aa) open reading

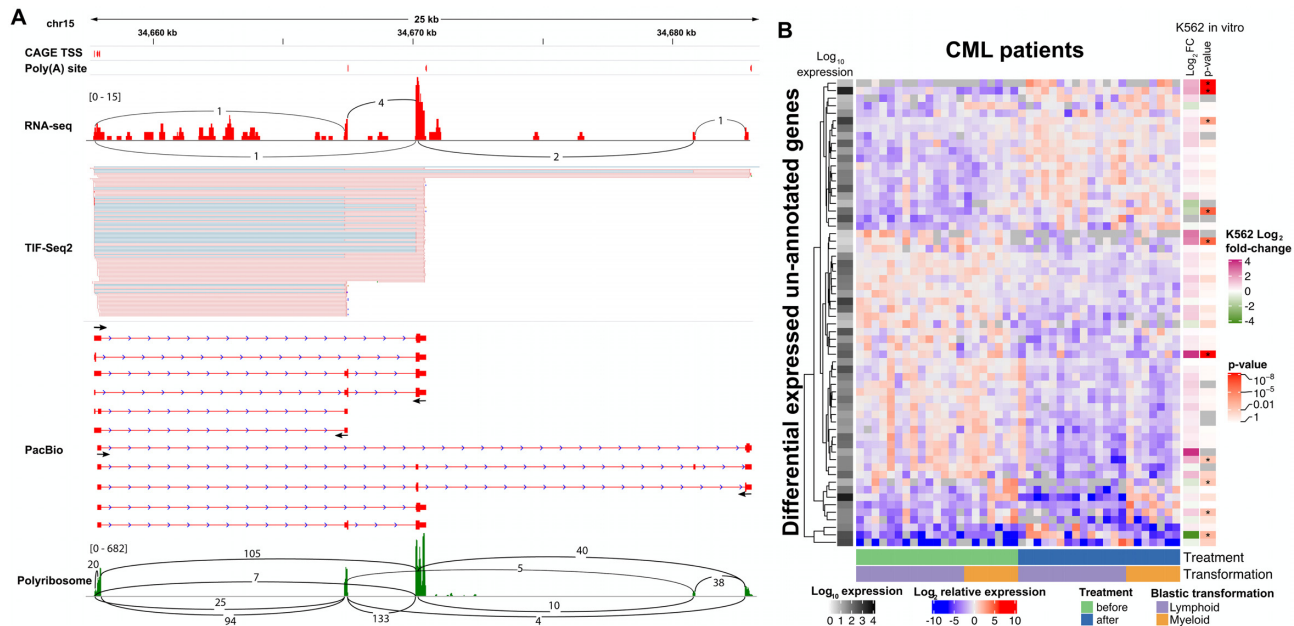
frames in 17 cases (Figure 2A and Supplementary Table S7), suggesting their coding potential.

After annotating those intergenic transcripts in K562, we investigated their potential relevance in chronic myeloid leukaemia (CML) patients. We used the TIF-Seq2-derived annotation to re-analyse RNA-Seq data from a cohort of CML patients responded to first-line TKI imatinib therapy, and focused on those who failed to achieve durable major molecular response and eventually developed blast crisis (27). We confirmed that 365 unannotated genes identified in K562 are also expressed in at least three CML patients with minimum one count per million reads. 59 of them responded to imatinib treatment, and 10 are differentially regulated between patient groups that developed myeloid blast crisis (MBC) or lymphoid blast crisis (LBC) (Figure 2B, adjusted  $P$ -value < 0.005). Interestingly transcripts down-regulated by imatinib treatment in K562 are upregulated in patients that stopped responding to treatment and underwent blast crisis. As an example, a K562-specific transcript on chr6:15555763–15559977 (putatively encoding the HTH domain of the Mos1 transposase) was significantly upregulated in patients who developed LBC after therapy, while it is significantly downregulated in K562 cell line after imatinib treatment (Supplementary Figure S13A–C). Independently of the potential relevance of the newly identified transcripts for disease progression, it is clear that they are present but excluded from most analysis due to incomplete transcriptome annotation and difficulty of its analysis.

### Disentangling overlapping transcription units allows unequivocal assignment of promoter proximal poly(A) sites

After showing the ability of TIF-Seq2 to better identify the boundaries of lowly expressed transcripts, we focused on the analysis of overlapping transcription units. We first asked how common was the overlap between neighbouring same-strand transcription units and their degree of overlap. Forty percent of overlapping TIF pairs present a high degree (over 90%) of overlap (top right corner in Figure 3A). This represents alternative transcription isoforms in same genes with slightly different TSSs or PASs in the first or last exons. About 50% of the TIF pairs represent the overlap of a longer isoform and its truncated short isoform, which starts from an intronic TSS or ends at an intronic PAS (top or right side in Figure 3A respectively). The rest of the overlapping TIF pairs correspond to tandem short isoforms originating from either the same gene loci or tandem transcripts originating from upstream genes. As partial overlaps between transcription units offer clear opportunities for regulatory crosstalk among adjacent genes (11,36), we decided to focus on those.

First, we focused on overlapping TIFs originating within 10 kb of each other. We can observe that upstream overlapping transcription units tend to use a poly(A) site located within the first 2 kb of the downstream transcription unit (x-axis in Figure 3B). And among the overlapping tandem TIFs originating in a 10 kb window from the downstream transcription unit, 60% of them arise within 1 kb upstream of the second TIFs (y-axis in Figure 3B). This suggests that in those regions, two RNA polymerases with different origins coexist, even with constitutive nu-



**Figure 2.** Unannotated lowly expressed intergenic transcripts can be detected in CML patients. (A) TIF-Seq2 can identify lowly expressed transcripts, for example, an unannotated gene on chr15: 34657737–34683026. TIF-Seq2 read pairs are labelled as pink lines (the forward strand). The fine blue lines in the TIF-Seq2 track represent splicing junctions in the first or last exons. CAGE TSS and poly(A) sites from public repositories (12,13) and an independent K562 RNA-Seq experiment (26) validates the expression of this gene. Primer designed for target long-read Pac Bio validation are marked by arrows. Polyribosome-associated RNA-Seq (28) track is in the bottom, with expression marked in green and splicing junction in black lines, suggests its coding potential. (B) expression of unannotated transcriptional features in CML patient (27). Each row represents a differentially expressed gene (adjusted  $P$ -value < 0.005) in CML patient. Each column represents a sample before or after drug treatment. Patients were classified as lymphoid blast crisis (LBC) or myeloid blast crisis (MBC) according to the types of their blastic transformations. Each cell represents  $\log_2$ -scale ratio of gene expression in each sample to the average expression of the gene across all samples. The average expression (in  $\log_{10}$  scale) of genes is depicted on the left. On the right side, we present the  $\log_2$ -scale fold change of those genes in an independent *in vitro* experiment exposing K562 cells to imatinib (26) and their significant levels of differential expression after treatment. Genes in K562 cells with adjusted  $P$ -value < 0.01 are labelled in asterisks.

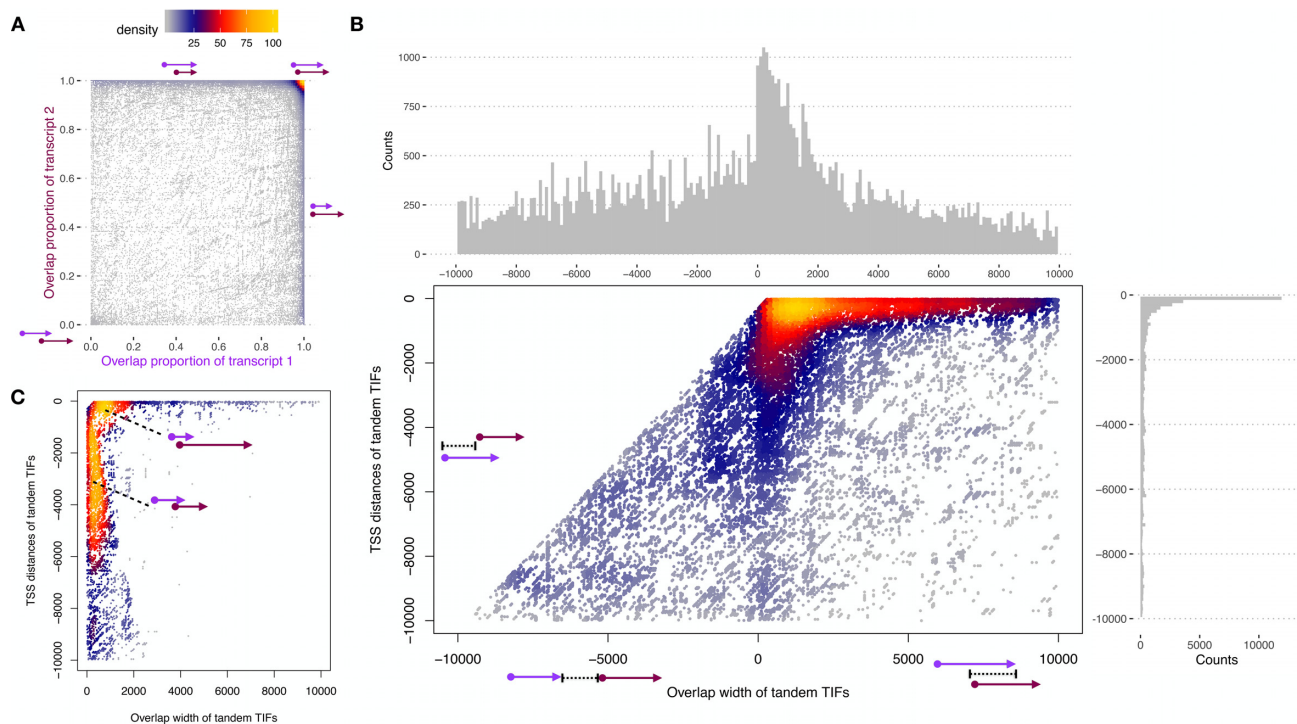
clear RNA degradation machinery. The first RNA polymerase originating from upstream promoter will terminate shortly after passing the start site of the second one, while the second RNA polymerase will proceed until the canonical poly(A) site. When investigating the case of low-degree tandem overlapping TIFs (<20% overlap in either TIF), we discovered two typical types of overlap (Figure 3C). One pattern (23%) is an upstream short TIF (<2 kb) originating from 1 kb promoter region of a downstream long TIF (as shown for gene GABPB1-AS1 in Supplementary Figure S14). Another type (30%) consists of two relatively longer TIFs that overlap at low degree. Polyadenylation cleavage sites of the upstream TIFs in both types typically occur within 2 kb downstream of the next TIF start sites, which matches the observation in all tandem overlapping TIFs (Figure 3B). In order to examine potential regulatory crosstalk among overlapping transcripts, we measured pairwise gene co-expression pattern in 21 CML patients RNA-Seq data (27). As RNA polymerases usually go further after poly(A) cleavage sites, we included the non-overlapping transcripts in the upstream 10 kb region. The TIF pairs are classified into different groups according to their overlap widths or the distance between two TIFs. Expression of TIFs are represented by its corresponding genes in the CML patients' data. In general, the overlapping transcripts shows a slightly less negative correlation than non-overlapping pairs (Supplementary Figure S15). However, no significant difference in co-expression was detected among groups.

Having an advantage of capturing both transcription boundaries in overlapping transcripts, TIF-Seq2 is able to pinpoint interdependence of TSSs and PASs across overlapping transcripts (Supplementary Figure S14). In addition, we observed a type of TSS–PAS combination from tandem neighbouring genes, thus generating fused genes caused by transcriptional read-through.

### Linking transcription boundaries facilitates the identification of read-through transcripts

An extreme case of overlapping transcription units, is the case of those where we identified novel combination between TSS–PAS from neighbouring genes. As those events suggest transcriptional read-through and the potential generation of fused genes, we decided to study them in detail. Transcriptionally fused genes are in general challenging to identify by conventional RNA-Seq approaches (37–39). Gene fusions play a key role in oncogenesis (40), and although most known fusion transcripts arise through chromosomal rearrangements, they can also arise via transcription-induced chimeras (e.g. read-through or alternative splicing) (37). Short-read RNA-Seq can detect read-through transcripts by using the small proportion of reads connecting two neighbouring genes (38). However, this becomes extremely challenging when the fusion event involves an intermediate ‘stepping stone’ exon not included in the annotations of the involved genes (e.g. LHX6-





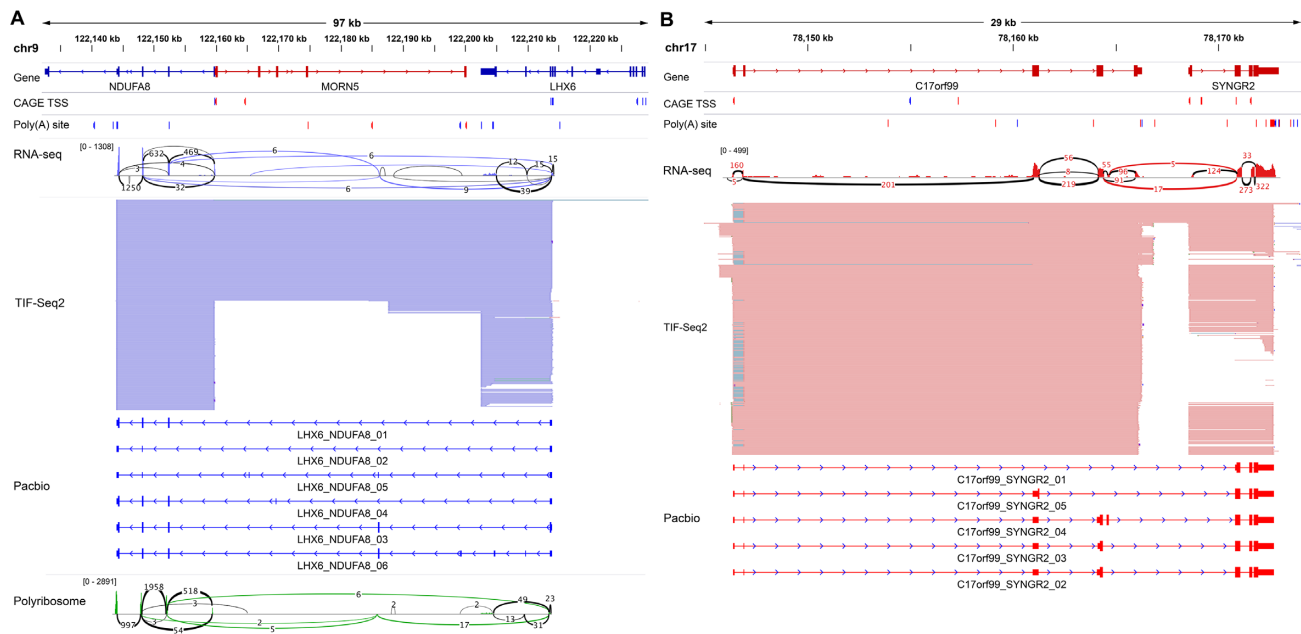
**Figure 3.** TIF-Seq2 reveals overlapping transcription units. (A) Pairwise analysis of the genomic overlap between TIFs. TIFs from the same gene loci present a high degree of overlap (over 90%, top-right corner), while TIFs from neighbouring tandem genes overlap at low degree (bottom-left corner). (B) Tandem TIFs captured by TIF-Seq2 reveal the path of RNA polymerase II (Pol II). Each dot in the heat scatter plot represents a pair of tandem TIFs within 10 kb. X axis shows overlap width between two tandem TIFs. In the case of no overlap, the values in X axis is the distance between two TIFs in the tandem TIF pairs. Y axis represents the distance of upstream TSS to the downstream TSS of tandem TIFs. Tandem TIFs typically overlap by less than 2 kb (histogram on the top). The majority of the first TIFs originated from within 1 kb upstream of the second TIFs. (C) Low degree overlapping tandem TIFs. Dots, x axis and y axis are as in (B).

NDUFA8.03, Figure 4A). Additionally, RNA-Seq does not allow to identify which TSS of the upstream gene is connected to a particular PAS of the downstream gene. On the contrary, all TIF-Seq2 reads have the potential to detect such fusion events (not only the small fraction covering a splicing between the annotated genes) and define their complete boundaries (TSS to PAS). We discovered 29 unannotated read-through transcripts in K562 cell, and once defined their boundaries, we were able to identify supporting splicing reads for all of them using RNA-Seq (26) (Supplementary Table S8). However, without the additional information provided by TIF-Seq2, those same RNA-Seq datasets were not sufficient to confidently classify them as read-through transcripts. Interestingly, during the preparation of this work, a few cases have been investigated in detail and independently reported as fusion transcripts in agreement with our findings (41,42).

To investigate up to what degree alternative splicing contributes to the appearance of those read-through transcripts, we further validated 10 candidates using the TIF-Seq2-derived boundaries and amplicon-based enrichment of full-length isoforms followed by long-read sequencing (as before, Supplementary Figure S12). This revealed an interleaved and complex transcriptome organization (Figure 4A, B and Supplementary Figure S16). In some cases, read-through transcripts connect the 5' regions of one gene with the coding sequence of the downstream one (e.g. C17orf99.SYNGR2.01, Figure 4B), suggesting a regula-

tory rewiring of the downstream gene (e.g. putative regulators of C17orf99 could regulate SYNGR2 expression), while in other cases, the read-through transcripts have the potential to encode fusion proteins (e.g. SPN\_QPRT\_04, Supplementary Figure S16). Our approach can detect even more complex scenarios, as is the case of transcripts initiating in the body of one gene but using the PAS of a downstream gene (e.g. LHX6-NDUFA8 in Figure 4A), or other complex cases where the read-through transcript connects both genes using an unannotated exon (e.g. LHX6-NDUFA8.03-06 or C17orf99.SYNGR2.04 in Figure 4).

To provide additional supporting evidence for the identified read-through transcripts, we investigated their potential to encode fusion proteins. We used RNA-Seq from polyribosome associated mRNAs (28) and confirmed that reads connecting the novel splicing sites predicted by our long-read experiments can be identified in 10 read-through genes (Figure 4A). This suggests that the identified transcripts associate to polyribosomes, and thus have the potential to encode fusion proteins. In addition, we investigated their presence in CML patients by re-analysing clinical RNA-Seq data (27). Even if we had to restrict our analysis to those few reads bridging the gene pairs, we were able to confirm their expression in patients (Supplementary Figure S17). Thus, by combining an improved annotation using TIF-Seq2 with available clinical short read dataset, we were able to identify read-through transcripts that would be ideal candidates for in-depth molecular characterization in disease models.



**Figure 4.** TIF-Seq2 facilitates the identification of read-through transcripts. (A) *LHX6-NDUFA8* read-through gene. (B) *C17orf99-SYNGR2* read-through gene. Annotated genes are listed on the top, followed by CAGE TSS and poly(A) site (PAS) (12,13) tags which are represented by red or blue bars (forward or reverse strands respectively). RNA-Seq validates the presence of splicing junctions (red or blue lines with supporting number of reads) linking two adjacent genes. TIF-Seq2 track (as in Figure 2A) shows the transcriptional fusion events between adjacent genes. PacBio Long-read sequencing of target transcripts validates the intergenic splicing events and dissect the transcription model of read-through genes. Polyribosome-associated RNA-Seq (28) data are labelled in green, with splicing junctions in green lines between two genes, showing the coding potential of two transcript isoforms of *LHX6-NDUFA8*.

## DISCUSSION AND CONCLUSION

Here we have presented an improved approach, TIF-Seq2, designed to link transcription boundaries in complex genomes. This approach is in agreement with previous maps of TSS and PAS, and can enrich current transcriptome studies by clarifying complex arrangement of overlapping transcripts. By focusing its sequencing power on TSS and PAS, it can easily define the complete boundaries of lncRNAs or other lowly expressed transcriptional features. This allows also to link genes with putative regulatory regions and to facilitate genetic manipulations (e.g. using CRISPRi or CRISPRa). As TIF-Seq2 is based on short-read sequencing, in principle it could be easily combined with traditional target enrichment approaches to link unannotated, but experimentally validated, TSS or PAS sites (i.e. using a unique probe targeting TSS or PAS of interest, without the need to enrich the intervening region). Therefore, it opens the door for the design of target enrichment strategies to allow their study with other sequencing approaches. By combining TIF-Seq2 information with RNA-Seq, we showed how an improved transcriptome annotation can refine our analysis of available clinical RNA-Seq datasets. This would be especially valuable to prioritize the transcriptional features more relevant for in-depth molecular validation. Application of TIF-Seq2 goes beyond the study of the human genome, and it could be particularly useful to facilitate the annotation of less studied genomes.

TIF-Seq2 is particularly useful to dissect those overlapping transcripts which pose a challenge to investigate by any other short-read RNA-Seq approaches. It can easily distinguish overlapping transcript units from each other and

quantify their expression, thus facilitating the interpretation of isoform-specific response in different conditions. By examining the interaction between overlapping transcripts, we observed that partially overlapping transcripts often use poly(A) sites within 2 kb of the TSS of the downstream genes. This suggests that during this window the cellular machinery is able to distinguish between short transcripts that should be terminated from those that will proceed to produce a full RNA. This is reminiscent of the mechanism previously described for PROMPTs and the transcription between closely spaced promoters (20,43). We showed that partially overlapping transcripts are formed by two main classes of upstream transcripts: those originating proximal to the downstream TSS (less than 1 kb) and a second-class involving transcript originating upstream. Those arrangements are difficult to study by RNA-Seq and CAGE and easily lead to the misassignment of reads corresponding to the upstream promoter to the downstream transcription unit. Finally, we identified read-through transcripts by linking the usage of particular TSS and PAS from neighbouring genes. This reveals the ability of read-through transcripts to putatively rewire their regulation (i.e. usage of alternative shared upstream promoters) and produce fusion proteins. We confirmed the identified fusion transcripts by targeted long-read sequencing and confirmed their putative coding potential by reanalysing ribosome profiling data. We showed that the identified transcripts can also be observed in clinical RNA-Seq datasets. Our work also points to the intrinsic limitation of RNA-Seq, which in some cases needs to be complemented by alternative approaches to resolve complex overlapping transcriptional structures. In conclu-

sion, we think that TIF-Seq2 has great potential to complement the current transcriptomic approaches, help dissect the overlapping transcriptome and thus fill in missing puzzle pieces to improve our understanding of transcription regulation.

## DATA AVAILABILITY

All TIF-Seq2, 3'T-fill sequencing and long-read sequencing files were deposited in the Gene Expression Omnibus under accession number GSE140912.

The source code for TIF-Seq2 data analysis is available on GitHub (<https://github.com/jingwen/TIFseq2>).

## SUPPLEMENTARY DATA

Supplementary Data are available at NAR Online.

## ACKNOWLEDGEMENTS

We thank Donal Barrett for technical assistance, Yuanyuan Xi for supporting computational analysis, and all members of Pelechano, Kutter and Friedländer laboratories for discussion. We thank Aino Järvelin and Judith Zaugg for early discussion. We would like to thank Dr Andreas Schreiber for his help with the access to clinical RNA-Seq datasets. The authors would like to acknowledge support from Science for Life Laboratory, the National Genomics Infrastructure (NGI) and Uppmax for providing assistance in massive parallel sequencing and computational infrastructure.

*Author contributions:* V.P., B.L., J.W., W.W. and L.S. conceived the study. B.L. and S.M. performed experiments and assisted with data interpretation. J.W., B.L., W.W. and V.P. designed data analysis and interpreted data. J.W. performed computational analysis. J.W., B.L. and V.P. drafted the original manuscript. All authors reviewed and edited the manuscript. V.P. supervised the whole project.

## FUNDING

Swedish Research Council [VR 2016-01842]; Wallenberg Academy Fellowship [KAW 2016.0123]; Swedish Foundations' Starting Grant (Ragnar Söderberg Foundation); Karolinska Institutet (SciLifeLab Fellowship, SFO and KI funds to V.P.); National Key R&D Program of China [2017YFC0908405]; National Natural Science Foundation of China [81870187 to W.W.]; US National Institutes of Health [NIH grant P01 HG000205]; Deutsche Forschungsgemeinschaft [1422/4-1]; European Research Council Advanced Investigator Grant (to L.M.S.). V.P. and W.W. acknowledge the support from a Joint China-Sweden mobility grant from STINT [CH2018-7750]; National Natural Science Foundation of China [81911530167]. Funding for open access charge: Vetenskapsrådet.

*Conflict of interest statement.* None declared.

## REFERENCES

1. Yao, P., Potdar, A.A., Arif, A., Ray, P.S., Mukhopadhyay, R., Willard, B., Xu, Y., Yan, J., Saidu, G.M. and Fox, P.L. (2012) Coding region polyadenylation generates a truncated tRNA synthetase that counters translation repression. *Cell*, **149**, 88–100.
2. Wei, W., Hennig, B.P., Wang, J., Zhang, Y., Piazza, I., Sanchez, Y.P., Chabbert, C.D., Adjalley, S.H., Steinmetz, L.M. and Pelechano, V. (2019) Chromatin-sensitive cryptic promoters putatively drive expression of alternative protein isoforms in yeast. *Genome Res.*, **29**, 1974–1984.
3. Tian, B. and Manley, J.L. (2017) Alternative polyadenylation of mRNA precursors. *Nat. Rev. Mol. Cell Biol.*, **18**, 18–30.
4. de Klerk, E. and 't Hoen, P.A.C. (2015) Alternative mRNA transcription, processing, and translation: insights from RNA sequencing. *Trends Genet.*, **31**, 128–139.
5. Uszczynska-Ratajczak, B., Lagarde, J., Frankish, A., Guigó, R. and Johnson, R. (2018) Towards a complete map of the human long non-coding RNA transcriptome. *Nat. Rev. Genet.*, **19**, 535–548.
6. Chabbert, C.D., Eberhart, T., Guccini, I., Krek, W. and Kovacs, W.J. (2019) Correction of gene model annotations improves isoform abundance estimates: the example of ketohexokinase (KHK). *F1000Research*, **7**, 1956.
7. Bertin, N., Mendez, M., Hasegawa, A., Lizio, M., Abugessaisa, I., Severin, J., Sakai-Ohno, M., Lassmann, T., Kasukawa, T., Kawaji, H. *et al.* (2017) Linking FANTOM5 CAGE peaks to annotations with CAGEscan. *Sci. Data*, **4**, 170147.
8. Hon, C.-C., Ramilowski, J.A., Harshbarger, J., Bertin, N., Rackham, O.J.L., Gough, J., Denisenko, E., Schmeier, S., Poulsen, T.M., Severin, J. *et al.* (2017) An atlas of human long non-coding RNAs with accurate 5' ends. *Nature*, **543**, 199–204.
9. Reyes, A. and Huber, W. (2017) Alternative start and termination sites of transcription drive most transcript isoform differences across human tissues. *Nucleic Acids Res.*, **46**, 582–592.
10. Proudfoot, N.J. (1986) Transcriptional interference and termination between duplicated  $\alpha$ -globin gene constructs suggests a novel mechanism for gene regulation. *Nature*, **322**, 562–565.
11. Van Werven, F.J., Neuert, G., Hendrick, N., Lardenois, A., Buratowski, S., Van Oudenaarden, A., Primig, M. and Amon, A. (2012) Transcription of two long noncoding RNAs mediates mating-type control of gametogenesis in budding yeast. *Cell*, **150**, 1170–1181.
12. FANTOM Consortium and the RIKEN PMI and CLST, D.G.T., Forrest, A.R.R., Kawaji, H., Rehli, M., Baillie, J.K., de Hoon, M.J.L., Haberle, V., Lassmann, T., Kulakovskiy, I.V., Lizio, M. *et al.* (2014) A promoter-level mammalian expression atlas. *Nature*, **507**, 462–470.
13. Gruber, A.J., Schmidt, R., Gruber, A.R., Martin, G., Ghosh, S., Belmadani, M., Keller, W. and Zavolan, M. (2016) A comprehensive analysis of 3' end sequencing data sets reveals novel polyadenylation signals and the repressive role of heterogeneous ribonucleoprotein C on cleavage and polyadenylation. *Genome Res.*, **26**, 1145–1159.
14. Byrne, A., Beaudin, A.E., Olsen, H.E., Jain, M., Cole, C., Palmer, T., DuBois, R.M., Forsberg, E.C., Akesson, M. and Vollmers, C. (2017) Nanopore long-read RNAseq reveals widespread transcriptional variation among the surface receptors of individual B cells. *Nat. Commun.*, **8**, 16027.
15. Volden, R., Palmer, T., Byrne, A., Cole, C., Schmitz, R.J., Green, R.E. and Vollmers, C. (2018) Improving nanopore read accuracy with the R2C2 method enables the sequencing of highly multiplexed full-length single-cell cDNA. *Proc. Natl. Acad. Sci. U.S.A.*, **115**, 9726–9731.
16. Wilkening, S., Pelechano, V., Järvelin, A.I., Tekkedil, M.M., Anders, S., Benes, V. and Steinmetz, L.M. (2013) An efficient method for genome-wide polyadenylation site mapping and RNA quantification. *Nucleic Acids Res.*, **41**, e65.
17. Martin, M. (2011) Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet journal*, **17**, 10.
18. Smith, T., Heger, A. and Sudbery, I. (2017) UMI-tools: modeling sequencing errors in unique molecular identifiers to improve quantification accuracy. *Genome Res.*, **27**, 491–499.
19. Dobin, A., Davis, C.A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., Batut, P., Chaisson, M. and Gingeras, T.R. (2013) STAR: ultrafast universal RNA-seq aligner. *Bioinformatics*, **29**, 15–21.
20. Chen, Y., Pai, A.A., Herudek, J., Lubas, M., Meola, N., Järvelin, A.I., Andersson, R., Pelechano, V., Steinmetz, L.M., Jensen, T.H. *et al.* (2016) Principles for RNA metabolism and alternative transcription initiation within closely spaced promoters. *Nat. Genet.*, **48**, 984–994.
21. Haberle, V., Forrest, A.R.R., Hayashizaki, Y., Carninci, P. and Lenhard, B. (2015) CAGER: precise TSS data retrieval and high-resolution promoterome mining for integrative analyses. *Nucleic Acids Res.*, **43**, e51.



22. Balwierz,P.J., Carninci,P., Daub,C.O., Kawai,J., Hayashizaki,Y., Van Belle,W., Beisel,C. and van Nimwegen,E. (2009) Methods for analyzing deep sequencing expression data: Constructing the human and mouse promoterome with deepCAGE data. *Genome Biol.*, **10**, R79.
23. Love,M.I., Huber,W. and Anders,S. (2014) Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol.*, **15**, 550.
24. Li,H. (2018) Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*, **34**, 3094–3100.
25. Quinlan,A.R. and Hall,I.M. (2010) BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics*, **26**, 841–842.
26. Gallipoli,P., Giotopoulos,G., Tzelepis,K., Costa,A.S.H., Vohra,S., Medina-Perez,P., Basheer,F., Marando,L., Lisio,L. Di, Dias,J.M.L. *et al.* (2018) Glutaminolysis is a metabolic dependency in FLT3ITD acute myeloid leukemia unmasked by FLT3 tyrosine kinase inhibition. *Blood*, **131**, 1639–1653.
27. Branford,S., Wang,P., Yeung,D.T., Thomson,D., Purins,A., Wadham,C., Shahrin,N.H., Marum,J.E., Nataren,N., Parker,W.T. *et al.* (2018) Integrative genomic analysis reveals cancer-associated mutations at diagnosis of CML in patients with high-risk disease. *Blood*, **132**, 948–961.
28. Ramirez,O., Kesarwani,A., Abhishek,G., Minella,A.C. and Pillai,M.M. (2016) Integrative analysis of RNA-Interactome and translome reveal functional targets of MSI2 in myeloid leukemia. *Blood*, **128**, 1881–1881.
29. Kim,D., Langmead,B. and Salzberg,S.L. (2015) HISAT: a fast spliced aligner with low memory requirements. *Nat. Methods*, **12**, 357–360.
30. Liao,Y., Smyth,G.K. and Shi,W. (2014) featureCounts: an efficient general purpose program for assigning sequence reads to genomic features. *Bioinformatics*, **30**, 923–930.
31. Pelechano,V., Wei,W. and Steinmetz,L.M. (2013) Extensive transcriptional heterogeneity revealed by isoform profiling. *Nature*, **497**, 127–31.
32. Matsumoto,K., Suzuki,A., Wakaguri,H., Sugano,S. and Suzuki,Y. (2014) Construction of mate pair full-length cDNAs libraries and characterization of transcriptional start sites and termination sites. *Nucleic Acids Res.*, **42**, e125.
33. Ruan,X. and Ruan,Y. (2012) Genome Wide Full-Length Transcript Analysis Using 5' and 3' Paired-End-Tag Next Generation Sequencing (RNA-PET). In: *Transcriptional Regulation. Methods in Molecular Biology (Methods and Protocols)*, vol 809. Springer, NY, pp. 535–562.
34. Pelechano,V., Wei,W., Jakob,P. and Steinmetz,L.M. (2014) Genome-wide identification of transcript start and end sites by transcript isoform sequencing. *Nat. Protoc.*, **9**, 1740–1759.
35. Lagarde,J., Uszczyńska-Ratajczak,B., Carbonell,S., Pérez-Lluch,S., Abad,A., Davis,C., Gingeras,T.R., Frankish,A., Harrow,J., Guigo,R. *et al.* (2017) High-throughput annotation of full-length long noncoding RNAs with capture long-read sequencing. *Nat. Genet.*, **49**, 1731–1740.
36. Pelechano,V. and Steinmetz,L.M. (2013) Gene regulation by antisense transcription. *Nat. Rev. Genet.*, **14**, 880–893.
37. Akiva,P., Toporik,A., Edelheit,S., Peretz,Y., Diber,A., Shemesh,R., Novik,A. and Sorek,R. (2006) Transcription-mediated gene fusion in the human genome. *Genome Res.*, **16**, 30–36.
38. Kumar,S., Razzaq,S.K., Vo,A.D., Gautam,M. and Li,H. (2016) Identifying fusion transcripts using next generation sequencing. *Wiley Interdiscip. Rev. RNA*, **7**, 811–823.
39. Hu,X., Wang,Q., Tang,M., Barthel,F., Amin,S., Yoshihara,K., Lang,F.M., Martinez-Ledesma,E., Lee,S.H., Zheng,S. *et al.* (2018) TumorFusions: an integrative resource for cancer-associated transcript fusions. *Nucleic Acids Res.*, **46**, D1144–D1149.
40. Mertens,F., Johansson,B., Fioretos,T. and Mitelman,F. (2015) The emerging complexity of gene fusions in cancer. *Nat. Rev. Cancer*, **15**, 371–381.
41. Weber,L., Maßberg,D., Becker,C., Altmüller,J., Ubrig,B., Bonatz,G., Wölk,G., Philippou,S., Tannapfel,A., Hatt,H. *et al.* (2018) Olfactory receptors as biomarkers in human breast carcinoma tissues. *Front. Oncol.*, **8**, 33.
42. Wu,P., Yang,S., Singh,S., Qin,F., Kumar,S., Wang,L., Ma,D. and Li,H. (2018) The landscape and implications of chimeric RNAs in cervical cancer. *EBioMedicine*, **37**, 158–167.
43. Ntini,E., Järvelin,A.I., Bornholdt,J., Chen,Y., Boyd,M., Jørgensen,M., Andersson,R., Hoof,I., Schein,A., Andersen,P.R. *et al.* (2013) Polyadenylation site-induced decay of upstream transcripts enforces promoter directionality. *Nat. Struct. Mol. Biol.*, **20**, 923–928.