

RESEARCH PAPER



Identification of RNA 3' ends and termination sites in *Haloferax volcanii*

Sarah J. Berkemer ^{a,b}, Lisa-Katharina Maier ^c, Fabian Amman ^{d,e}, Stephan H. Bernhart ^{a,f}, Julia Wörtz^c, Pascal Märkle^c, Friedhelm Pfeiffer ^g, Peter F. Stadler ^{a,h,i,j,b,d,k,l}, and Anita Marchfelder ^c

^aBioinformatics Group, Department of Computer Science - and Interdisciplinary Center for Bioinformatics, University of Leipzig, Leipzig, Germany; ^bMax Planck Institute for Mathematics in the Sciences, Leipzig, Germany; ^cBiology II, Ulm University, Ulm, Germany; ^dInstitute for Theoretical Chemistry, University of Vienna, Vienna, Austria; ^eDivision of Cell and Developmental Biology, Medical University Vienna, Vienna, Austria; ^fTranscriptome Bioinformatics, Interdisciplinary Center for Bioinformatics, Leipzig University, Leipzig, Germany; ^gComputational Biology Group, Max Planck Institute of Biochemistry, Martinsried, Germany; ^hFacultad de Ciencias, Universidad Nacional de Colombia, Bogotá, Colombia; ⁱCenter for RNA in Technology and Health, University Copenhagen, Frederiksberg C, Denmark; ^jSanta Fe Institute, Santa Fe, NM, USA; ^kGerman Centre for Integrative Biodiversity Research (iDiv), Halle, Jena and Leipzig, Germany; ^lCompetence Center for Scalable Data Services and Solutions, and Leipzig, Research Center for Civilization Diseases, University Leipzig, Leipzig, Germany

ABSTRACT

Archaeal genomes are densely packed; thus, correct transcription termination is an important factor for orchestrated gene expression. A systematic analysis of RNA 3' termini, to identify transcription termination sites (TTS) using RNAseq data has hitherto only been performed in two archaea, *Methanosarcina mazei* and *Sulfolobus acidocaldarius*. In this study, only regions directly downstream of annotated genes were analysed, and thus, only part of the genome had been investigated. Here, we developed a novel algorithm (Internal Enrichment-Peak Calling) that allows an unbiased, genome-wide identification of RNA 3' termini independent of annotation. In an RNA fraction enriched for primary transcripts by terminator exonuclease (TEX) treatment we identified 1,543 RNA 3' termini. Approximately half of these were located in intergenic regions, and the remainder were found in coding regions. A strong sequence signature consistent with known termination events at intergenic loci indicates a clear enrichment for native TTS among them. Using these data we determined distinct putative termination motifs for intergenic (a T stretch) and coding regions (AGATC). *In vivo* reporter gene tests of selected TTS confirmed termination at these sites, which exemplify the different motifs. For several genes, more than one termination site was detected, resulting in transcripts with different lengths of the 3' untranslated region (3' UTR).

ARTICLE HISTORY

Received 4 December 2019
Revised 23 January 2020
Accepted 24 January 2020

KEYWORDS

Transcription termination; archaea; Haloarchaea; *Haloferax volcanii*; RNAseq; RNA 3' ends; 3' UTR

Introduction

Archaeal RNA synthesis is generally considered to be more closely related to transcription in eukaryotes than to bacterial transcription. The archaeal RNA polymerase is similar to the eukaryotic RNA polymerase II, and the basal promoter elements in Archaea are similar to their eukaryotic pendants (TATA box and BRE); general transcription factors TBP (TATA binding protein), TFB (transcription factor B) and TFE (transcription factor E) resemble the eukaryotic proteins TBP, TFIIB and TFEII, respectively (for a review see: Fouqueau et al. [1]). Transcriptional regulators, however, seem to be more similar to those in bacteria [2]. Thus, the archaeal transcription machinery consists of a mixture of bacterial-like and eukaryotic-like components. Whereas some data have been reported on transcription initiation and elongation in archaea, very little is known about transcription termination. Controlled transcription termination is important to avoid aberrant RNA molecules and to help with RNA polymerase recycling. Generally, the genes in archaeal chromosomes are densely packed, so that proper termination is also important to prevent transcription from continuing into downstream genes. The process of

transcription termination is not trivial because the very stable transcription elongation complex must be destabilized and dissociated during termination. In bacteria, two major classes of termination signals have been described: intrinsic termination and factor-dependent termination [3,4]. Intrinsic termination occurs either at a stretch of Ts or at hairpin structures that fold in the newly synthesized RNA; both trigger dissociation of the elongation complex. Factor-dependent termination occurs upon interaction with a specific protein such as the bacterial termination factor Rho [5]. Protein factor-assisted termination is especially important in regions where strong selective pressure on the DNA sequence does not allow encoding of intrinsic termination signals. This may be the case when termination must occur in the coding region of a downstream gene. In eukaryotes, several RNA polymerases synthesize the different RNA classes, and these polymerases have different modes of termination. RNA polymerase I requires protein factors for termination [6], whereas RNA polymerase III terminates efficiently and precisely at a stretch of Ts [7]. Termination of RNA polymerase II is quite complex, involving modification of the RNAP as well as interactions with additional

CONTACT Peter F. Stadler  studla@bioinf.uni-leipzig.de  Bioinformatics Group, Department of Computer Science - and Interdisciplinary Center for Bioinformatics, University of Leipzig, Härtelstraße 16-18, Leipzig D-04107, Germany; Anita Marchfelder  anita.marchfelder@uni-ulm.de  Biology II, Ulm University, Ulm 89069, Germany

 Supplemental data for this article can be accessed [here](#).

© 2020 The Author(s). Published by Informa UK Limited, trading as Taylor & Francis Group. This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

protein factors, and seems to be coupled to co-transcriptional RNA processing [8]. Whereas archaeal transcription initiation resembles eukaryotic RNA polymerase II initiation [9], transcription elongation and termination seem to be more similar to the eukaryotic RNA polymerase III pathway, which is independent of RNA secondary structures and protein co-factors [10].

Compared to the determination of transcription start sites, the identification of termination sites is more complex. Termination is often leaky and encompasses several consecutive sites, and degradation by exonucleases renders the 3' ends heterogeneous and less clear. Also, original transcription termination sites are difficult to distinguish from RNA 3' termini that are due to RNA processing. Data reported on archaeal termination so far show that intrinsic termination occurs at a run of Ts [10–14] and is potentially also influenced by secondary structure elements [10,14]. Factor-dependent termination was predicted based on the results of an *in vivo* reporter assay in the archaeon *Thermococcus kodakarensis* [15]. A recent study confirmed this hypothesis, reporting the discovery of the first archaeal termination factor [16].

Recently, RNA 3' ends for two archaea (*Sulfolobus acidocaldarius* and *Methanosarcina mazei*) were investigated systematically using RNAseq data [17]. All identified RNA 3' ends were considered to reflect transcription termination sites (TTS) and TTS were found for 707 and 641 transcriptional units in *S. acidocaldarius* and *M. mazei*, respectively. For more than a third of the genes analysed, multiple consecutive terminators were identified, resulting in 3' untranslated regions (UTRs) with different lengths [17]. In some cases, the terminator of an early gene in an operon was shown to be located in the downstream gene, allowing gene-specific regulation within an operon. The study also revealed lineage-specific features of termination for both archaea, confirming the requirement that more data from different archaeal organisms are required to learn more about transcription termination in this domain. However, the applied algorithm analysed only regions directly downstream of annotated genes, and thus only a part of the genome was considered. Total cellular RNA was used for the analyses and no attempt was made to distinguish between RNA 3' ends originating from transcription termination and RNA processing.

Here, we describe the identification of RNA 3' ends of the halophilic model archaeon *Haloferax volcanii* and the subsequent determination of termination motifs. *H. volcanii* has been used for a plethora of biological studies [18,19], including the determination of nucleosome coverage [20] and a genome-wide identification of TSS [21]. *H. volcanii* requires high salt concentrations for optimal growth, and due to the high intracellular salt concentrations, RNA-protein interactions -including modes of transcription termination- may differ from those in mesophilic archaea.

We used a newly developed algorithm to identify RNA 3' ends from RNAseq data genome-wide in an unbiased manner, independent of annotation and on the basis of reads enriched for primary transcripts. Applying this algorithm to RNAseq data from a terminator exonuclease (TEX) treated library, we found 1,543 RNA 3' ends for the *Haloferax* genome. Subsequent analysis of the respective sites revealed a strong sequence signature consistent with the current archaeal

termination model, indicating a strong enrichment of native transcription termination sites (TTS) among the newly identified RNA 3' ends. Therefore, identified RNA 3' ends were considered putative TTS for successive analysis with respect to primary and/or secondary structure motifs as termination signals and 3' UTR characterization. Selected termination motifs were confirmed using an *in vivo* reporter gene system.

Results

Identification of RNA 3' ends downstream of annotated regions

To identify RNA 3' ends downstream of annotated regions in *H. volcanii* we applied a self-implemented version of the recently published method from Dar et al. [17], which will be referred to as the Dar-Sorek-Method (DSM). RNA was isolated from *H. volcanii* cells (from three biological replicates) and cDNA libraries were generated to allow determination of RNA 3' ends. Libraries were made such that the original RNA 3' end is tagged and can be identified in the resulting sequence. Libraries were subjected to paired-end next-generation sequencing (NGS), resulting in 49 million reads for each library on average. 89–98% of the reads obtained mapped to the *Haloferax* genome (Supplementary Table 6). Since total RNA was used for library preparation the identified RNA 3' ends result from either transcription termination or processing. Mapped reads were analysed with DSM, which identifies RNA 3' ends in a defined region downstream of annotated genes at the position with the highest coverage of mapped read ends. The length of the downstream region is determined by the average insert length of corresponding paired-end reads (Fig. 1).

In our dataset, the median length of the analysed region was 126 bp. The resulting 3' UTRs were mostly shorter than 100 nucleotides, with a median length of 58 nucleotides (Supplementary Figure 1). The length restriction given in the DSM approach might be too strict for genes with long 3' UTRs. In addition, the method cannot determine RNA 3' ends independent of an annotation (see also the paragraph below 'Comparison between DSM and IE-PC algorithms'). Using DSM, we identified 3,155 RNA 3' ends for the complete *Haloferax* genome of which 85% were in intergenic and the remainder in coding regions (Supplementary Table 3). A typical RNA 3' end is shown in Supplementary Figure 2 for the *pilA2* gene (HVO_2062).

Development of a novel algorithm, 'internal enrichment-peak calling', as a tool to identify RNA 3' ends genome-wide

DSM analysis only includes sequences downstream of annotated genes and, thus, only a fraction of the genome [17]. To overcome this restriction, we developed a novel approach to interpret the RNAseq data obtained which we termed Internal Enrichment-Peak Calling (IE-PC). To this end, we utilized the nature of the fragmentation process in the course of library preparation, which leads to an enrichment of natural 3' fragment ends over 5' ends of fragments that are generated via random

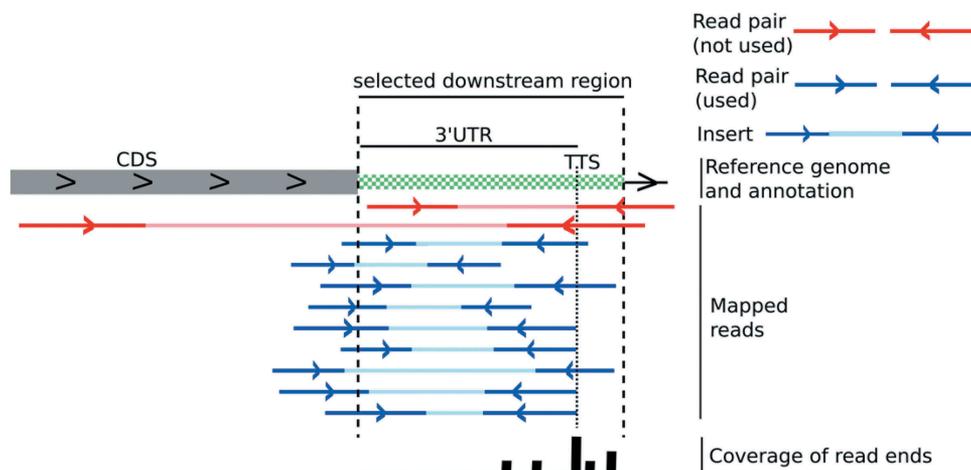


Figure 1. Principle of the Dar-Sorek-Method. As a first step in the DSM method, read pairs with an insert overlapping an annotated region are selected (red and blue lines) [17]. Inserts that do not overlap or have a length of more than 500 nucleotides are discarded (red read pairs in the Figure). From the selected read pairs (blue), the average length over all inserts is calculated as described in Dar et al [17]. This value is used to determine the length of the selected downstream region (green white region in the Figure). The coverage of all read ends in this region is retrieved (bottom line), and the position with the highest coverage is identified as the TTS.

fragmentation during library preparation. Since we used the latter as the normalization background for the former, we called this method internal enrichment (IE) (Fig. 2). This approach requires mate-pair sequencing data, from which a set of read-pairs is selected that contains the same original first-strand cDNA 5' ends (which reflects the RNA 3' end) (Fig. 2). The corresponding 3' ends of the mate-pair reads are then evaluated. Sequencing the cDNAs generated from these first-strand cDNA fragments in paired-end mode preserved information about fragment ends. After mapping the read-pairs to the reference genome, the hallmark of an original RNA 3' end was its high coverage of read ends, with its associated mate ends originating from a multitude of genomic sites. It is highly unlikely that independent clones end at an identical position (Fig. 2). Multiple fragments with a heterogeneous 5' end but a common 3' end thus are indicative of true natural RNA 3' ends.

To gain further specificity, putative 3' ends were subsequently censored if the site in question was associated with a sudden decrease in read coverage, indicative of a true RNA 3' end. This was assessed by a peak calling approach (termed PC), considering the absolute and relative number of fragments ending at the respective position (Fig. 3). The two methods (IE and PC) were sequentially run on the data, and only sites that were found by both independent approaches, IE as well as PC, were considered to be bona-fide RNA 3' ends. We allowed a maximal distance of 10 nucleotides when computing the intersection of IE and PC. The advantage of this algorithm is that it is independent of genome annotation and thus analyses the complete genome sequence rather than a restricted annotation-dependent region allowing the identification of all RNA 3' ends of a genome.

Comparison between DSM and IE-PC algorithms

Since our newly established method IE-PC works independent of any genome annotation, it is able to determine RNA 3' ends covering the complete genome. Comparison of signals obtained with both methods (DSM and IE-PC) can only be

done with the regions that are included in the DSM analysis. The original DSM analysis was based on regions with a median length of 126 nucleotides downstream of an annotated 3' end, with the RNA 3' end located at the position with the highest coverage.

We compared individual RNA 3' ends found with DSM and/or IE-PC in downstream regions defined by the DSM analysis in more detail. Examples are shown in Figs. 4 and 5, Table 1 lists the number of RNA 3' ends that were identified by both methods (allowing for a position discrepancy of up to 10 nucleotides) as well as RNA 3' ends that appeared in only one of the data sets. Altogether, we found 1,664 RNA 3' ends being present in both data sets (Table 1).

Comparing the DSM data with IE-PC results as well as with given coverage data, we see that chosen downstream regions in DSM were frequently too short.

This is also clear when comparing the median 3' UTR length determined by both methods. While DSM found a median 3' UTR length¹ of 58 nucleotides, IE-PC determined it to be 97 nucleotides. A comparison of RNA 3' ends found with DSM and IE-PC for the HVO_1876s gene is shown as an example (Fig. 4).

The lower panel (DSM) in Fig. 4 shows the results obtained with DSM, coverage is shown as the 3' end coverage of corresponding reads, the RNA 3' end position for the transcript is also shown. The upper panel (IE-PC) shows IE-PC results, with RNA 3' end location and read coverage above. The drop in coverage values is clearly visible, matching the RNA 3' end location as identified by IE-PC. Corresponding read end coverage of the DSM data set was present at the same position. However, in the DSM analysis, every gene had a downstream region with an individual length, since for each gene the average insert length of the corresponding reads was taken into account. The length for the downstream region of HVO_1876s, was calculated with only 64 nucleotides, and therefore the high read end coverage was beyond the region selected for analysis. As the algorithm is bound to report an RNA 3' end (there is always a position with highest coverage in the selected region), DSM reported a false positive in this case. For the genes *trmY* and *tRNA^{Pro}*, both methods

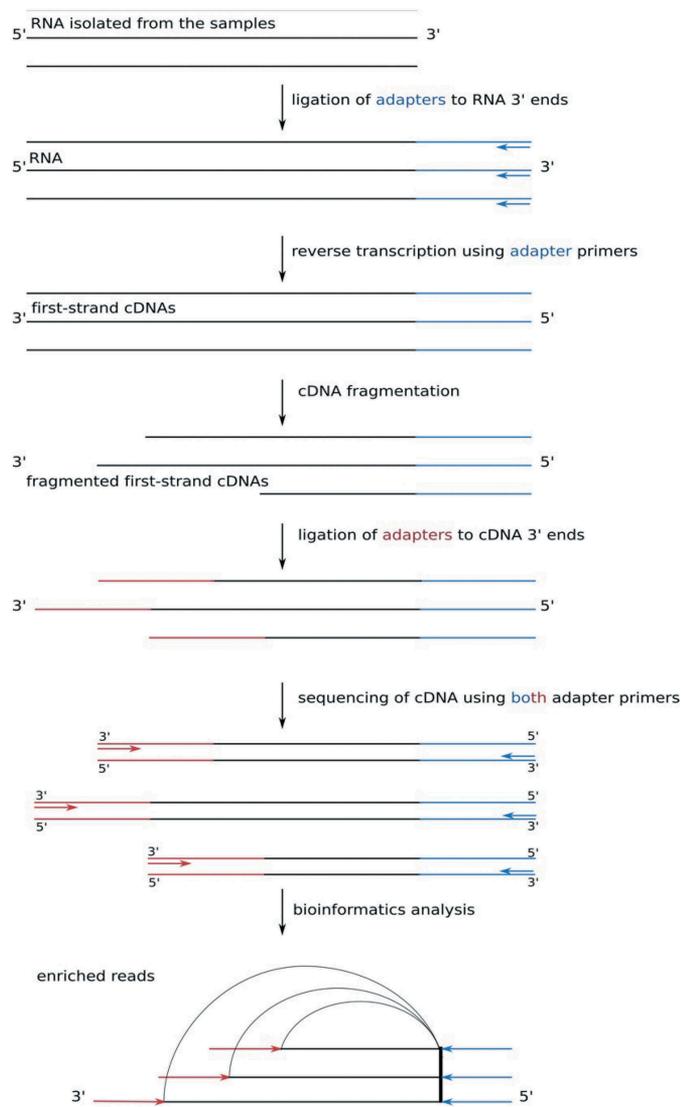


Figure 2. Principle of internal enrichment. After RNA isolation, adapter primers were ligated to the RNA 3' end and the RNA was reverse transcribed. First-strand single-stranded cDNA was fragmented prior to the addition of the adapter primer at the cDNA 3' end (for details, see materials and methods and supplementary methods). The break points were considered to be random, leading to an enrichment of original RNA 3' ends over 5' fragmentation ends. Sequencing the cDNAs generated from these first-strand cDNA fragments in paired-end mode preserved information about fragment ends, even if they were longer than the read length. After mapping the read-pairs to the reference genome, the hallmark of an original 3' end was its high coverage of read ends, with its associated mate ends originating from a multitude of genomic sites.

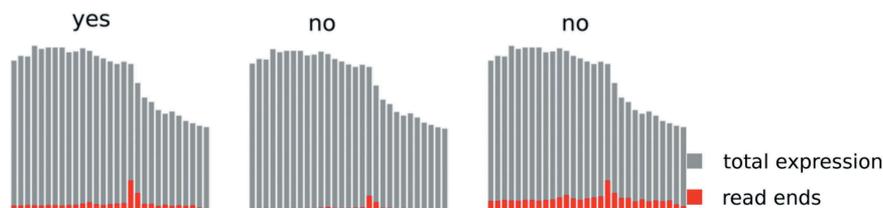


Figure 3. Principle of the applied peak calling procedure. In a sliding window approach, positions where the number of read stops (red) exceeded the mean number of read stops over the whole window by a z-score above the threshold 2 were treated as potential endings (left), while potential peaks of the same height but below a z-score of 2 were discarded (right). In addition, if the number of read stops was below a 10% threshold relative to the total coverage (grey), the peak was discarded (centre).

identified one identical RNA 3' end position and two different RNA 3' ends positions (Fig. 5). The IE-PC analysis identified two RNA 3' ends downstream of the $tRNA^{Pro}$ gene (upper panel). The corresponding coverage values for DSM (lower panel) showed similar signals, but due to its specific algorithm, DSM assigned one RNA 3' end downstream of the *trmY* gene (HVO_1989) and one downstream of the $tRNA^{Pro}$ gene (lower panel).

Taken together comparison between the DSM approach and the IE-PC algorithm clearly shows that using the IE-PC approach yields improved and more comprehensive data.

Novel approach for the identification of transcription termination sites

To determine RNA 3' ends the DSM approach used reads from a total cellular RNA fraction that contains RNA 3' ends derived from transcription termination as well as 3' ends derived from processing. Thus the 3' ends identified by DSM are not all TTS but also processing sites (PS). A similar problem exists for the determination of original transcription 5' ends, where a well established method for the reliable identification of transcription start sites (TSS) has been developed, termed differential RNAseq (dRNAseq). Here, to enrich primary transcripts with original 5' ends an RNA sample is treated with terminator exonuclease (+TEX) which removes RNAs with a 5'-monophosphate. Primary transcripts are newly synthesized RNA molecules that have not been processed at their 5' end and also have a higher probability to contain the original 3' terminus. The majority of bacterial ribonucleases prefer substrates with 5'-monophosphate ends, some even have a specific sensor domain for the 5'-monophosphate [22], thus primary transcripts are less prone to processing.

Therefore -similar to the approach for start site determination- we used a +TEX library to enrich original termination ends. To that end we treated a cellular RNA fraction with 5' terminator exonuclease (TEX) to enrich primary transcripts and thereby original termination ends. After cDNA library generation from the TEX treated RNA, NGS was performed, resulting in an average of 40 million reads for each of the three libraries (Supplementary Table 6). Reads obtained were analysed with our newly established algorithm as described below.

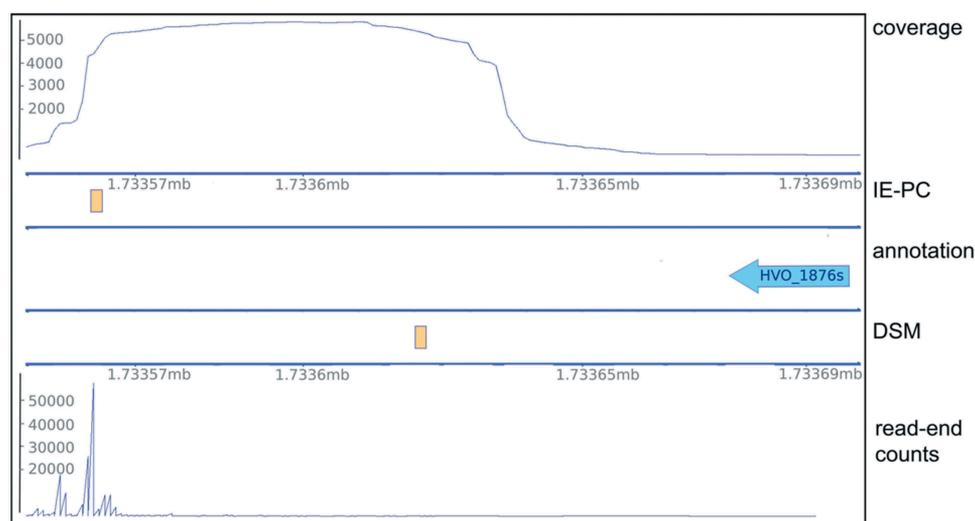


Figure 4. TTS comparison for DSM and IE-PC analyses for the transcript HVO_1876s. The data in the lower panel (DSM) show the DSM results with TTS and corresponding read end counts from the DSM data set. The upper panel (IE-PC) shows the TTS locations determined by IE-PC and the total coverages, corresponding to the -TEX data set. Since sequencing starts at the 3' end, coverage starts at the 3' end and runs continuously for 75 bp due to the read length. The annotation and genome coordinates are shown in the middle (coordinates given in Mb). TTS are shown as orange rectangles.

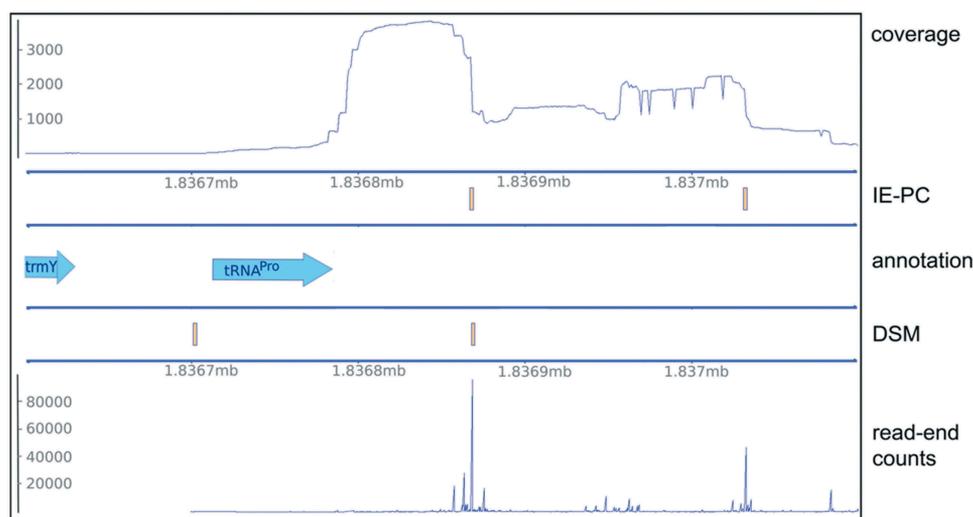


Figure 5. TTS comparison of DSM and IE-PC for the $tRNA^{Pro}$ gene. In the lower panel, termination sites and corresponding read end coverages from the DSM data set are shown, assigning a TTS downstream of the *trmY* gene and another one downstream of the $tRNA^{Pro}$ gene (shown as orange rectangles). The upper panel shows the TTS assignment determined by IE-PC and the total coverages, corresponding to the -TEX data set, identifying two TTS downstream of the $tRNA^{Pro}$ gene. Since sequencing starts at the 3' end, coverage starts at the 3' end and runs continuously for 75 bp due to the read length. The secondary TTS of the $tRNA^{Pro}$ gene was not reported by the DSM algorithm as this algorithm systematically reports only a single TTS for each annotated gene.

Table 1. Comparison of RNA 3' ends found with DSM and IE-PC in downstream regions covered by DSM. The column 'IE-PC only' lists all RNA 3' ends identified by IE-PC but not with those identified by DSM method. RNA 3' ends identified with IE-PC that were also found with DSM are listed in the column 'overlapping'. Sites that were detected in the defined region only by DSM and that did not overlap with IE-PC sites are listed in the column 'DSM only'. The IE-PC data shown here are the ones calculated on the basis of the -TEX data set, since the DSM data are also based on the -TEX data.

Chromosome	IE-PC only	Overlapping	DSM only
Main	3,184	1,296	1,128
pHV1	231	53	39
pHV3	394	97	125
pHV4	811	218	199
Total	4,620	1,664	1,491

The Haloferax genome contains 1,543 transcription termination sites

To identify transcription termination sites we applied the IE-PC algorithm to the data from the TEX treated samples consisting of enriched original transcription termination sites, and identified 1,543 putative TTS² (Table 2). Supplementary Table 1 lists all TTS detected.

Clustering of termination signals

Inspection of the TTS obtained revealed closely spaced TTS, that were 11 to 150 nucleotides apart. These closely spaced

Table 2. TTS identified with IE-PC. TTS identified are present in coding regions as well as in intergenic regions.

	TTS
Intergenic	807
Coding	736
Total	1,543

TTS were subclassified into first TTS (TTS₁) and secondary TTS (TTS_s): a first TTS is located directly downstream of a 3' gene end on the same strand (Fig. 6A, B); a secondary TTS is located downstream of another TTS, with no other features (like TSS or 3' gene end) in between, as shown in Fig. 6C. We found in total 1,056 first TTS and 487 secondary TTS. We also found very closely spaced TTS that were less than 10 nucleotides apart these were reported as only a single TTS (the one with the highest coverage).

Distinct motifs for transcription termination in coding and intergenic regions

From the 1,543 transcription termination sites found, slightly more than half of the sites were found in intergenic regions (807 TTS, 52%) and the remainder in coding regions (736 TTS, 48%) (Table 2). Detailed information for each TTS can be found in Supplementary Table 1.

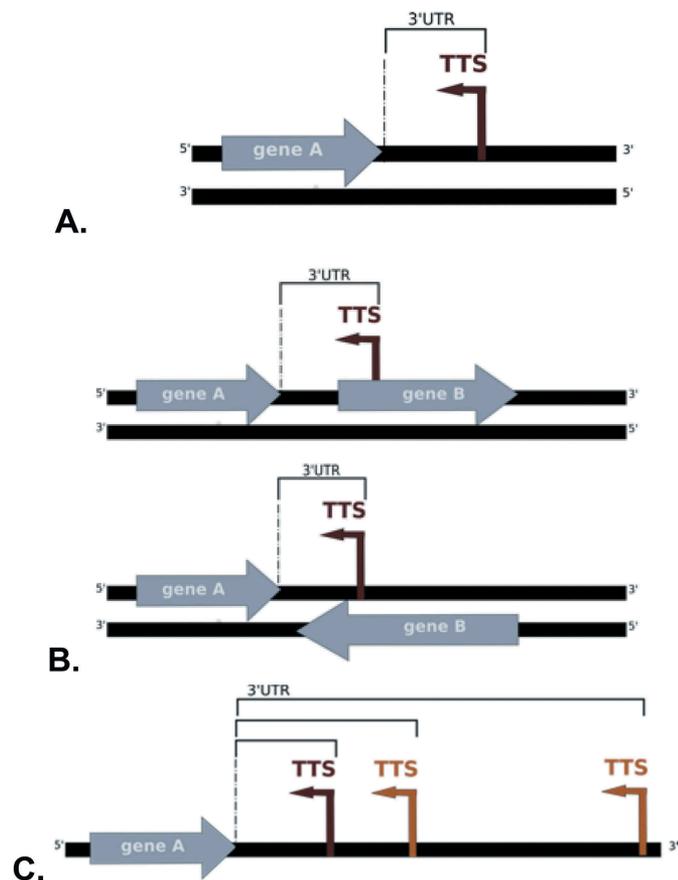


Figure 6. Location of TTS. A. TTS is located in an intergenic region. B. Location of the TTS in an annotated gene that is located on the sense or antisense strand. The length of the UTR was in all cases measured as the distance between the TTS and the 3' end of the upstream annotated gene on the same strand. C. A first (dark brown) and two secondary TTS (light brown) are shown.

Analysis of the regions up- and downstream of the TTS were performed separately for TTS located in coding and intergenic regions (Fig. 7), and in both, an increase in hybridization energy at the TTS similar to the increase identified in the TTS set obtained with DSM was found (Supplementary Figure 3). The pattern of nucleotide enrichment for sites located in intergenic regions showed that Ts were prevalent at the TTS (Fig. 7A).

We next analysed sequences 15 nucleotides up- and five nucleotides downstream of the 807 intergenic termination sites for common sequence motifs. For 748 sites, we found similarities in the sequences, such as a conserved dinucleotide TC as part of the motif as well as a stretch of T's of variable length upstream of the termination site (Fig. 8A, Supplementary Table 2).

Sequences 15 nucleotides up- and five nucleotides downstream of the 736 TTS in coding regions were likewise investigated for common sequence motifs (Fig. 8B). The prominent C residue at every third position is typical for coding regions in *Haloferax*. The third codon position combines an enrichment for GC (due to the GC-rich genome) and pyrimidines.

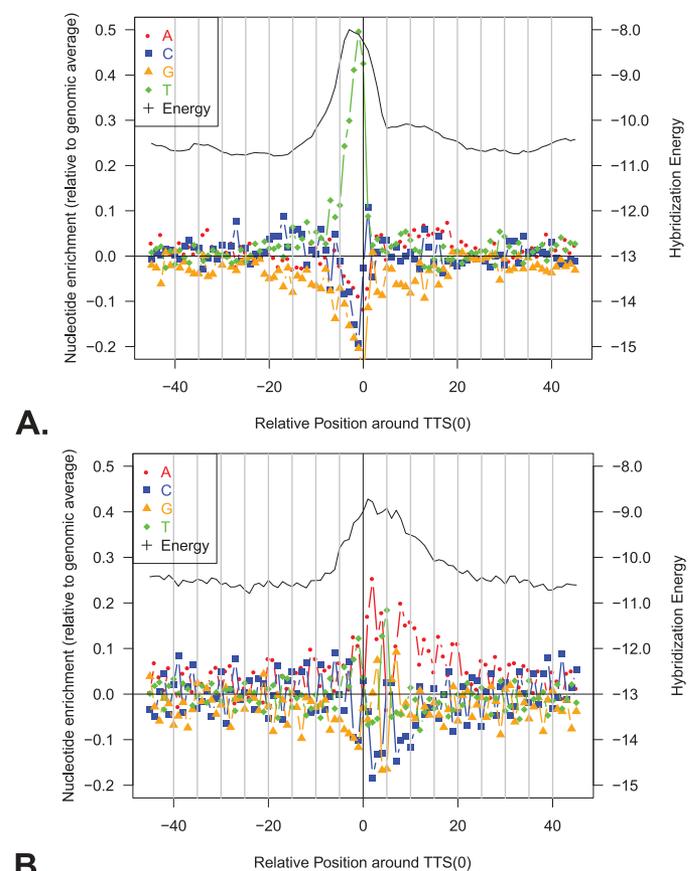


Figure 7. Analysis of up- and downstream regions for TTS identified with IE-PC. A. Intergenic and B. coding regions were investigated. Forty-five nucleotides up- and downstream of the termination site were analysed for (1) nucleotide enrichment at each position (left y-axis) and (2) the hybridization energy (right y-axis). x-axis: nucleotide position (upstream -, downstream +). The colour scheme for the four nucleotides is shown at the upper left, and the energy data are shown with a black line. Hybridization energies were calculated based on the binding energies between the DNA template and RNA in the area behind the RNA polymerase. The regular pattern in panel B is due to the statistical characteristics of coding regions in GC-rich organisms [23,24].

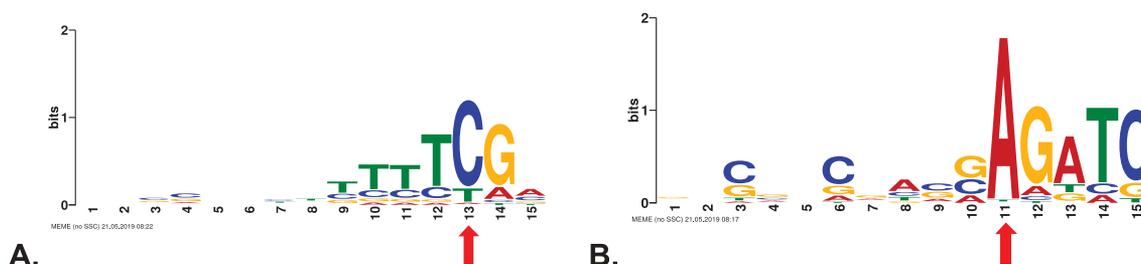


Figure 8. Enriched sequence motifs around the TTS identified with IE-PC. A. Sequence motif close to TTS located in intergenic regions. The TTS is located at position 13 (red arrow). B. Sequence motif close to TTS located in coding regions. The TTS is located at position 11 (red arrow). Motifs were detected using MEME [48].

For 456 TTS in coding regions, we found the downstream motif AGATC (Fig. 8B).

Taken together, we could identify distinct termination motifs that were specific for intergenic and for coding regions.

Secondary structures can act as termination signals

To identify potential secondary structure motifs, we conducted a search for accessible and inaccessible regions around the TTS using RNAplfold [25] (Vienna RNApackage [26,27]). We plotted accessibilities against a background distribution of shuffled dinucleotides. However, no clear signals were found to indicate significantly increased or decreased accessibility. In a second attempt to identify secondary structures, we applied graphclust 2.0 [28] within Galaxy [29] to sequences 100 nucleotides upstream of all TTS (Supplementary Figure 4). Graphclust is a tool that clusters input sequences based on their secondary structure(s). It will cut the sequences based on a window size parameter and align and fold the sequences into secondary structures using RNAalifold of the ViennaRNA package [26,27]. Using this approach we found hairpin structures upstream of 503 of the TTS (up to 10 nucleotides distance), an example is shown in Fig. 9.

Thus, in some cases secondary structures are present that might influence transcription termination.

Experimental confirmation of selected termination signals

We selected four TTS identified with our new IE-PC approach in intergenic and coding regions to test their termination activity with an *in vivo* reporter gene assay. Termination regions (the TTS and up- and downstream regions) were cloned into a fragment of the reporter gene β -galactosidase (sequences are listed in Supplementary Table 5). Termination activity was monitored with northern blot analyses (Fig. 10, Supplementary Table 8). If termination occurred in the inserted fragment, the RNA was shorter than that in the control construct (Fig. 10). The T stretch identified here and in earlier studies with archaea terminated efficiently at the expected site with 95% termination and some read-through, showing that this assay worked well (Fig. 10B, Supplementary Table 8). Next, we tested two TTS (A1 and A2) found in coding regions with the downstream motif AGATC. Both terminated at the expected site, confirming the activity of this newly identified downstream motif. However, termination is not as efficient as for the T stretch motif (11% for A1 and 27% for A2) and thus allows considerable read-through (Fig. 10C, Supplementary Table 8). Furthermore, a hairpin motif found upstream of TTS in coding regions (S12, shown in Fig. 9) was tested in the assay (Fig. 10D, Supplementary Table 8). Again, termination was detected at the expected site, 59% of the transcripts are terminated at the structural motif. To exclude that RNAs were processed at the sequence or structure

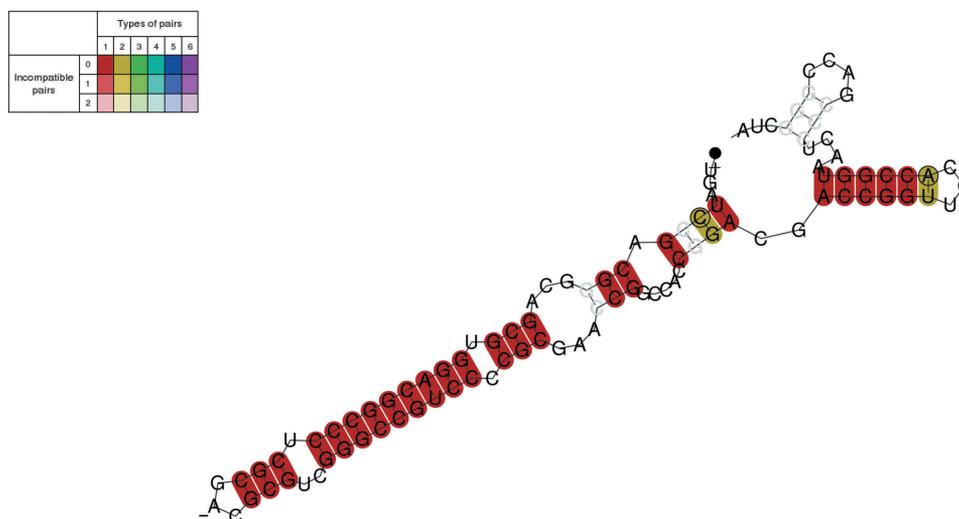


Figure 9. Hairpin structure found upstream of TTS. The secondary structure was plotted with RNAalifold, and its 5' end is denoted by a small dot at the end of the line. The colour coding is at the top left; whereas darker colours show compatible pairs, and the number of types of pairs shows how many types are found at this position thereby indicating sequence conservation. Dark red colours indicate base pairs that are compatible and conserved.

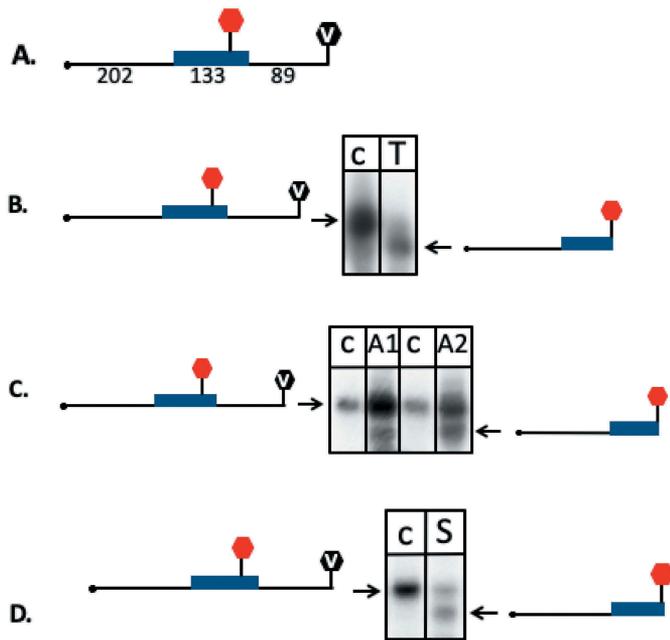


Figure 10. *In vivo* reporter gene test of termination sites. Termination sites were cloned into a reporter gene (the gene for β -galactosidase) construct to confirm their termination activity. RNA was isolated from strains transformed with reporter gene plasmids, separated by size and subsequently transferred to a nylon membrane. Membranes were hybridized with a probe against the β -galactosidase mRNA. Active termination sites generated shorter RNA molecules than the control sequences, indicated by an arrow. Experiments were carried out in triplicate. A. The construct used is shown schematically. A full length transcript terminating at the vector termination site is 424 nucleotides long. The insert containing the termination site is 133 nucleotides long (the insert is shown in blue, the location of the termination site in the insert is indicated with a red hexagon), the transcribed plasmid sequences are 202 nucleotides (upstream of the insert) and 89 nucleotides (downstream of the insert) long. The termination site in the plasmid is indicated with a V in the black hexagonal. B. The T₅C termination sequence was inserted into the reporter gene. Such a T₅C sequence can be found for example in the intergenic sequence downstream of the *fts2* gene. Transcripts terminating at the T₅C motif are 321 nucleotides long. The T stretch terminated efficiently (indicated with an arrow) with 95% termination at the T stretch and 5% termination at the plasmid encoded terminator (lane T), lane c: control insert without a termination site. C. Two sequences with the downstream motif AGATC (A1 and A2) were inserted into the reporter gene. Both are located in coding regions; A1 is located in the gene HVO_2724, A2 is located in the gene HVO_0455. Termination at these motifs results in 309 nucleotide long RNAs. Both sequences terminated, however considerable read-through is present (termination at A1: 11% and at A2: 27%). Lane A1: AGATC termination site A1, lane A2: AGATC termination site A2. D. A sequence with the potential to fold into a hairpin structure (S12) was inserted into the reporter gene. The structure S12 is located in a coding region (gene HVO_A0420). Termination at this motif results in a 309 nucleotide long RNA. S12 terminated, but considerable read-through was observed (termination at S12: 59%). Lane S: sequence S12.

motifs, northern blots were hybridized with a probe against the downstream fragment. If transcription would not terminate at the insert sequence but promote to the final terminator present in the vector and the transcript would be subsequently processed by an endonuclease, the downstream processing product would be picked up with this hybridization. However, no additional RNA fragments were detected (Supplementary Figure 6).

Analysis of 3' UTR length

To determine 3' UTR lengths, the distance between the first TTS and the 3' end of the preceding annotated gene was determined. The first TTS had a median 3' UTR length of

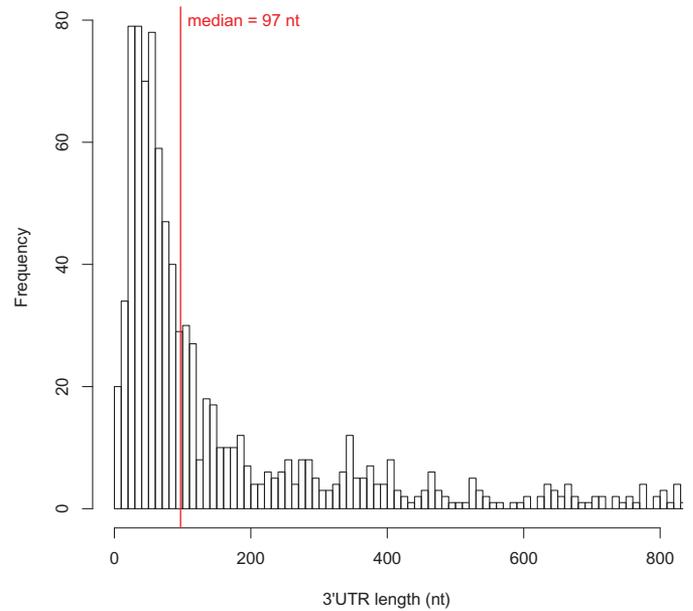


Figure 11. Histogram of 3' UTR lengths. The 3' UTR lengths for all first TTS are shown as calculated from the IE-PC data on the basis of the +TEX reads. The histogram shows the frequency of 3' UTR lengths for the region up to 800 nucleotides in length.

97 nucleotides (Fig. 11). The median 3' UTR length for all TTS (first and secondary) is 189 nucleotides. To confirm that the 3' UTR regions were part of the same transcripts as the upstream coding regions, we generated transcriptome data. To that end we performed next-generation sequencing of a cDNA library generated from total RNA for RNAseq. An average of 42 million reads were obtained for three independent cDNA libraries, and with these data, we confirmed that the 3' UTR regions were part of the same transcripts as the upstream coding regions, since we found continuous RNAseq reads over the complete putative 3' UTRs (Supplementary Table 4).

Expression of individual genes in operons

Earlier analyses of TSS data from *H. volcanii* revealed that the transcription start site for a gene in an operon can be located in an upstream gene that is part of the same operon [21]. Similarly, the TTS data reported here show that transcription can be terminated in a downstream gene in an operon. Together, this allows fine-tuning of the expression of individual genes in an operon and thus challenges the simplistic view that all genes of an operon are co-transcribed under all conditions. This phenomenon was confirmed in a previous study experimentally, for instance for the dicistronic operon encoding the proteins Lsm and L37e (HVO_2723 and HVO_2722). Here, a promoter for expression of the *rpl37e* gene is located in the upstream *lsm* gene (Supplementary Figure 5) [30]. In addition, the upstream *lsm* gene can be terminated within the downstream *rpl37e* gene. A dicistronic transcript and a *rpl37e* mRNA were identified by northern blot analysis [30]. A similar independent transcription of genes that are part of the same operon has been described in *Sulfolobus islandicus* [31]. To analyse the extend of the expression regulation of individual genes in an

operon we determined the operons of *H. volcanii* (Supplementary Table 9) and compared the transcript reads of the first and the last gene in each operon (Supplementary Figure 7). The majority of genes in an operon has the same amount of transcripts for the first and the last gene however with some clear exceptions that are listed in Supplementary Table 10 and 11. Operons with internal TTS show a slightly stronger tendency to have such outliers.

Discussion

The novel IE-PC approach allows comprehensive, genome-wide identification of transcription termination sites

Applying the new IE-PC algorithm on +TEX RNA sequences, we were able to determine the TTS genome-wide for the archaeon *H. volcanii*. The dRNAseq method using a +TEX library is an established method to confidently identify transcription start sites by enriching primary transcripts. However, treatment with TEX does not completely remove 5'-monophosphate RNAs, thus only an enrichment of primary transcripts is achieved and the +TEX data set contains to some extent still 5' processed RNA. Nevertheless, the method is standard and accepted for TSS identification.

Currently very little is known about archaeal RNases and their preferred substrates, but bacterial ribonucleases clearly prefer substrates with a 5'-monophosphate. Thus an RNA with an 5'-triphosphate is not the preferred substrate for these RNases. Therefore, we used here a +TEX library to enrich transcription termination sites. Altogether, we found 1,543 TTS in the +TEX data set with our new algorithm, representing the first unbiased, truly genome-wide approach.

Comparison of both methods for TTS identification, IE-PC and DSM, showed that IE-PC could identify TTS genome-wide, whereas DSM was restricted to regions downstream of annotated genes. Furthermore, algorithmic details of DSM (e.g. the requirement that each annotated gene is directly followed by a TTS; UTR length restriction by using average insert length) caused false TTS identifications in some cases.

Motifs for termination

Amongst all TTS, two general sequence motifs for termination were identified. One was predominant in intergenic regions, and the other one in coding regions. Both types of termination motifs (T stretch and AGATC) were confirmed using *in vivo* assays. Motifs in intergenic regions consisted of a stretch of Ts, while termination occurred at a C residue. This motif is similar to those found by *in vitro* studies in other archaea [10,12–14,32] and for *S. acidocaldarius* and *M. mazei* using the DSM algorithm [17].

A specific motif for coding regions

The presence of TTS in coding regions, terminating transcription of an upstream located gene, has been observed previously in archaea including the RNASeq analysis by Dar et al. [17,33]. Our novel IE-PC algorithm allowed to identify enough TTS in

coding regions to be able to identify a novel termination motif, AGATC. The *in vivo* termination test confirmed that AGATC is used as termination motif but revealed that it is not very efficient (11–27%). This lower termination efficiency is important for TTS in coding regions since it allows full length transcription of the gene in which the termination site is located. A termination event as strong as T₅C would be detrimental for the expression of the terminator containing gene. In addition, it would be very difficult to place such a T stretch in a coding region since the coding region is governed by the requirement to contain the specific triplets for the encoded protein, and especially a G/C rich genome like the one from *Haloferax* can rarely contain a T stretch in a coding region. Thus, to be able to terminate in coding regions other signals have to be used and these motifs should allow for sufficient read-through.

The prevailing motif we have found in coding regions was AGATC, located downstream of the termination site. The DNA sequence downstream of the TTS has been shown to influence termination efficiency in bacteria and eukaryotes [13]. Furthermore, it has been reported that the archaeal RNA polymerase interacts with downstream duplex DNA [34]. Thus, it is entirely possible that a downstream termination motif exists. It might also be possible that a protein binds to the downstream motif, acting as a minor roadblock for the RNA polymerase and thereby inducing pausing and in some cases subsequently terminating transcription. Future research will show the mechanistic details of this process.

Secondary structures involved in termination

The search for secondary structure motifs located close to the TTS was successful for a subset of the TTS and revealed hairpin structures upstream of the TTS, suggesting that these structures might influence transcription termination. One of the identified structures was tested in the *in vivo* termination experiments and indeed terminated transcription, however not as efficiently as the T stretch. A certain amount of folding potential upstream of TTS was also found in *M. mazei* but was considered not to be essential, whereas no structures were identified near TTS in *S. acidocaldarius* [17]. Possibly, termination signals in different archaeal phyla are more variable.

Different 3' UTR lengths are generated by termination at several TTS

With our genome-wide TTS determination approach, we can present the first comprehensive determination of 3' UTR lengths in *H. volcanii*. The median 3' UTR length for first TTS was 97 nucleotides, which is longer than the median 3' UTR length recently reported for two archaea (with 55 nucleotides for *S. acidocaldarius* and 85 nucleotides for *M. mazei*) and the bacterium *B. subtilis* (40 nucleotides) [17]. The presence of transcripts with long UTRs was confirmed by RNAseq data.

In many cases, we found several TTS downstream of one gene, showing that mRNAs with different 3' UTR lengths are generated. The lengths of the UTRs as well as the fact that some genes have several 3' UTRs provide a good platform for interactions with regulatory molecules, such as sRNAs. This is

similar to the situation in eukaryotes, in which termination at different sites generates different 3' UTRs, that can interact with different regulatory proteins or RNAs such as miRNAs. Different UTR lengths were also found in *M. mazei* and *S. acidocaldarius* [17].

Transcription produces antisense RNA

If a gene is located on the opposite strand of a TTS (Fig. 6B, lower panel), an antisense RNA against this oppositely encoded gene is generated. Of the 1,543 TTS detected, 14.4% (222 TTS) of them were located in a gene on the opposite strand, thus resulting in antisense RNAs. Up to 75% of all genes were found to be associated with antisense RNAs in bacteria [35–37]. Studies with archaea painted a similar picture [21,38]. In *S. acidocaldarius*, 301 gene pairs with convergent orientations were investigated, and in 52% of them, the terminator of one gene was located in the coding region of a gene on the opposite strand [17] (as in Fig. 6B lower panel); however, in *M. mazei*, only 8% of the convergent genes showed such an overlap. The potential functions and the physiological relevance of these antisense transcripts must be uncovered in the future.

Regulation of genes in operons

In *H. volcanii* transcription start sites [21] as well as termination sites (this work) were found in genes that are part of the same operon showing that genes in an operon can be individually regulated and expressed. A similar individual regulation of operonic genes has been shown in *S. islandicus* [31], suggesting that this is potentially a general mode of regulation in archaea.

Conclusion

Taken together, we showed that our new IE-PC approach used with a +TEX data set is well-suited to identifying the complete set of TTS for a genome without any a priori limitations on the search space due to genome annotation. We confirmed the T stretch termination motif detected 30 years ago, but we also identified an additional new motif specific for coding regions as well as a hairpin structure for termination. The presence of multiple 3' UTRs for a gene provides a platform for regulatory mechanisms similar to those described in eukaryotic systems.

Materials and methods

Haloferax volcanii culture conditions and RNA extraction

H. volcanii strains H119 ($\Delta leuB$, $\Delta pyrE2$, $\Delta trpA$) [39] and HV55 ($\Delta leuB$, $\Delta pyrE2$, $\Delta trpA$, $\Delta bgaH$) (this work, see below) were grown aerobically at 45°C in Hv-YPC or Hv-Ca medium [39] to an OD₆₀₀ of 0.8–0.9 (for a list of strains used see Supplementary Table 5). Total RNA was isolated from three biological replicates using TRIzol (ThermoFisher scientific). RNA fractions were sent to vertis (vertis Biotechnologie AG,

Martinsried, Germany) for further treatment, cDNA library preparation and high throughput sequencing. RNA preparations were made for each of the three different RNAseq approaches: (1) To obtain an RNA fraction enriched in primary transcripts RNA was treated with terminator exonuclease (this fraction was termed +TEX), (2) RNA from one preparation was left untreated representing the complete RNA pool (this fraction was termed -TEX). RNA preparations (1) and (2) were performed such that the original RNA 3' end was maintained (for details see Supplementary methods). For RNA preparation (3) RNA was isolated for obtaining transcriptome data (this fraction was termed RNAseq data). In all sequencing approaches primers were designed such that the identification of the RNA strand was possible. Details for library construction and sequencing are reported in Supplementary methods.

E. coli culture conditions

E. coli strains DH5 α (Invitrogen) and GM121 [40] were grown aerobically at 37°C in 2YT medium.

Read mapping

Raw reads were adapter clipped and quality trimmed using cutadapt version 1.10 [41] based on fastqc version 0.11.4 [42] quality control reports. Reads were then mapped with sege-mehl (version 0.2.0) [43,44]. We used -A 94 to require higher accuracy in order to account for prokaryote mapping instead of mapping eukaryotic genomes. Mapped reads were afterwards processed using samtools version 1.3 [45]. In order to calculate genome coverage and intersection of data sets, we used bedtools (bedtools v2.26.0) [46].

DSM analysis

Based on the Dar-Sorek-Method (DSM) [17,47], RNA 3' ends (putative TTS) were retrieved from the mapping data. We used samples without TEX treatment for consistency with the DSM analyses and we only used read pairs and at least a coverage of 4 for every valid position. As described in Dar et al. [17,47], for each annotated region in the genome, all inserts overlapping with the annotated region were collected. Then, for each position downstream of a gene, all collected reads that end at this position were summarized, excluding inserts longer than 500 bp. The average insert length of all the corresponding read pairs was defined as the length of the target downstream region. For all read pairs where the insert overlapped an annotated region, the position with highest read-end coverage inside the target downstream region of an annotated 3' end was reported as a transcription termination site (TTS) (Table 1 and Fig. 1). Collected TTS with DSM are listed in Supplementary Table 3. The used implementation of this method is available at Bioinformatics Leipzig (<http://www.bioinf.uni-leipzig.de/publications/supplements/18-059>).

Internal enrichment algorithm

To detect sites with a significant enrichment of sequenced and mapped fragment ends a sound background without enrichment is desired. In the current setting, we used the intrinsic properties of a paired-end sequencing run to directly deduce the following information (Fig. 2). We also took advantage of the fact, that the cDNA library preparation applies a fragmentation step after ligation of the 3' end primer. Since each fragment which results from an individual fragmentation event (in contrast to PCR duplicated fragments) is very unlikely to have the exact same length, truly enriched sites can be expected to be associated with sequenced fragments, all ending at the respective site but starting at different positions. Therefore, the more different mates (mapping to different positions) are associated with the different reads ending at a particular site the higher the enrichment of read end signal at that particular position can be considered. To capture this, we calculate for each position i a score S as

$$S_i = \frac{C_i}{\left(\prod_{j=1}^n C_j\right)^{1/n}} \forall_j \exists R(i \circ j)$$

Thereby, C_i denotes the number of fragment ends at position i , C_j the number of fragment starts at position j , and $R(i \circ j)$ all position i, j which are associated via at least one read-mate-pair R . To get an expected background distribution of these scores, we again use the nature of paired-end reads. Since we expect only fragment ends, in contrast to fragment starts, to be enriched, we can use the distribution of the reciprocally defined scores for the fragment start S_j as a background distribution. Based on the distributions, native fragment end scores can be assigned an empirical p-value, evaluating its likelihood to occur by chance alone. The above detailed software, named Internal-Enrichment (IE) was implemented in Perl and is available at Bioinformatics Leipzig (<http://www.bioinf.uni-leipzig.de/publications/supplements/18-059>). We ran IE using all position showing signals for starts and ends (-mr 1), we omitted read fragments with length more than 100 nt (-mf 100) to reduce outliers, and the geometric mean as a method to calculate a score for the given signals (-mode GeomMean). We only included signals that were present in all three replica such that each showed a maximal empirical p-value of 0.05. The geometric mean of all empirical p-values for one position from all replica is used as the final score for the given position.

Peak calling

A complementary approach to find RNA 3' ends is identifying peaks in read stops (Fig. 3). In order to find these peaks, we first computed the strand specific read coverage at every position in the genome, which we used as a background.

We then used a sliding window approach with a window size of 150 nt and an overlap of 50 nt to find positions in the respective windows with a significantly higher number of read stops than the rest of the window. This was done by computing the mean number of stops as well as the standard deviation by window and then calculating the z-score for the

number of read ends at every position of the window. Subsequently, we required a minimum number of 5 reads stopping, at least 10% of all reads covering position $i-1$ must end at position i as well as a minimum z-score of +2 at a site to report it as a putative RNA 3' end. As a consequence of the overlapping of the windows, a site is evaluated multiple times and must fulfil the criteria in at least one of the contexts evaluated.

Enrichment of primary transcripts

To be able to differentiate between 3' end processing sites and termination sites we isolated RNA and enriched this RNA for primary transcripts by digestion with the terminator exonuclease (TEX) (for details see Supplementary methods). Primary transcripts contain 5' triphosphates and are not digested by TEX. Bacterial ribonucleases preferentially use 5'-monophosphate RNA as substrates [22], thus the RNA fraction enriched for 5-triphosphate RNAs should be enriched in unprocessed RNAs.

Terminator identification

We report results for the sample enriched in primary transcripts (+TEX), in which IE and PC were both computed for the +TEX sample, followed by intersection analysis, resulting in the identification of TTS. In most cases, both methods (IE and PC) reported the same sites. For a few cases, more than one position was reported within a distance of 10 nt. In this case, the position with the highest coverage was chosen as TTS.

More information is available at Bioinformatics Leipzig (<http://www.bioinf.uni-leipzig.de/publications/supplements/18-059>).

Terminator sequence and structure analysis

We calculated percentages of nucleotide enrichment and hybridization energies for regions of 45 nt upstream and downstream of the TTS. Hybridization energies are calculated as the energy that stabilizes the RNA-DNA hybrid during transcription [25], using RNAplfold, a program part of the ViennaRNA package [27] (version 2.4.9) with special energy parameters for RNA-DNA hybrids. The lower (the more negative) the hybridization energy, the more stable is the RNA-DNA hybrid.

Motif and structure analysis for 3' UTR sequences

Motifs were detected using MEME [48] (version 5.0.1) by scanning all sequences from 15 nt upstream to 5 nt downstream of the TTS (see Fig. 8 for an example). Structure search was conducted using graphclust 2.0 [28] within Galaxy [29]. Input sequences were sequences of 100 nt length upstream to the TTS. Graphclust was used with default parameter and additionally a window size of 110 (such that our sequences fit in one window to avoid duplicated sequences), a bitscore of 15 for the results of cmscan, an upper threshold of 50 clusters and 20 top sequences in each alignment for the visualization. The application of graphclust resulted in 37 clusters providing

covariance models (CMs) for each. Covariance models are probabilistic models that are created based on given sequence and structure motifs and can be used to search for sequence and structure homologies. To scan sequences for these given motifs, we applied cmscan (a part of the infernal program suite [49], version 1.1.1) on our input sequences after calibrating the CMs (using cmcalibrate, part of the infernal program suite) and additionally on all the transcripts in order to get a background model. Given cmscan results, we filtered the secondary structures detected in the TTS-related sequences such that the structures were at most 10 nt upstream from the TTS.

Analysis of 3' UTR regions

We confirmed 3' UTR regions by checking if coding region and 3' UTR are completely covered by RNAseq reads. If so, this would be indicative of uninterrupted transcription until the assigned TTS. Out of 1,543 3' UTR sequences, 1,286 show a continuous coverage within RNAseq reads (82.2%) (Supplementary Table 4).

Code availability

The code is available at: <http://www.bioinf.uni-leipzig.de/publications/supplements/18-059> (see also data availability statement below).

Reporter gene investigations of identified TTS

Construction of *bgaH* deletion mutant HV55

To allow the use of a plasmid carrying the β -galactosidase reporter gene (a fusion of β -galactosidase genes from *Haloferax alicante* and *Haloferax volcanii*) (*bgaHa*) [50,51], the *H. volcanii* β -galactosidase gene (*bgaH*) was deleted in strain H119. H119 was transformed with the *bgaH* deletion plasmid pTA617 (for a list of plasmids used see Supplementary Table 5) [30,52] and grown in Hv-Ca medium with tryptophan (final concentration 0.25 mM) to generate the deletion strain [53]. Homozygous knock-out clones were verified via southern-blot using *SalI* digested genomic DNA and probe *bgaHaDO* (primers: *bgaHKODO*-for and *bgaHKODO*-rev (for a list of primers used see Supplementary Table 5); template genomic DNA of *H. volcanii*). Probe labelling and detection of the blot were carried out using the DIG-DNA labelling mix and detection reagents (Anti-Digoxigenin-AP) (Roche) according to the manufacturer's protocol. The resulting deletion strain was termed HV55.

Construction of pTA231-termtest constructs

To construct the terminator-test plasmid a *p.syn* expression cassette (Anice Sabag-Daigle and Charles J. Daniels, in preparation) synthesized by lifetechnologiesTM (ThermoFisher Scientific) was introduced via *NotI/EcoRI* into pTA231 [39], resulting in pTA231psyn (for sequences see Supplementary Table 5). This plasmid was subsequently cured of the *BamHI* site by digestion with *BamHI*, blunting by Klenow fragment and religation. Then a C-terminal fragment of the reporter gene *bgaHa* was inserted at the *NdeI* site. After *NdeI* digestion, the vector was treated with Pfu polymerase to fill-in the

NdeI site. The inserted *bgaHa* fragment was generated by PCR using primers *bgaHatermifw* and *bgaHatermirev* and pTA599 as template [52]. The correct orientation of the inserted fragment was confirmed by sequencing, the resulting plasmid was termed pTA231-termtest. The candidate terminator sequences were inserted into the *BamHI* site that is present in the newly inserted fragment. Terminator fragments were generated with PCR using primers as listed in Supplementary Table 5 and templates pTA599 (for the control construct) or genomic DNA of *H. volcanii* (for TTS-A1, TTS-A2, TTS-S12 and TTS-ftsZ). HV55 cells were transformed with pTA231termtest constructs, and grown in Hv-Ca medium with uracil (final concentration 0.45 mM). Cells were grown to an OD₆₅₀ of 0.6–1.0 and RNA was isolated using NucleoZOL (Machery-Nagel) according to manufacturer's instructions.

Northern-blot analysis of HV55xpTA231-termtest RNA

To analyse the RNA levels, total RNA was isolated from *H. volcanii* as described above. Ten μ g RNA was separated on a 1.5% agarose gel and subsequently transferred to a nylon membrane (Biodyne[®] A, PALL). After UV-crosslinking the membrane was hybridized with a radioactively labelled probe against the 5'-part of the *bgaHa* mRNA to detect termination events. The probe was generated by PCR (primers *bgaHatermi fw* and *TermiVectorrev*; template pTA599) and the purified PCR fragment was labelled using α -³²P-dCTP and the random primed DNA labelling kit DECAprimeTMII (Invitrogen). Experiments were done in triplicate. For detection of a potential processing fragment downstream of the termination site, a radioactively labelled probe against the 3'-part of the *bgaHa* mRNA was used. The probe was generated and labelled as described above using primers *BetaHinten1* and *BetaHinten2* and pTA231-termtestA1 as template. For quantification of northern blot signals, membranes were exposed to phosphorimaging plates (FujiFilm) and the resulting signals detected using a Typhoon imager (GE). Analysis of three replicates was carried out using the ImageQuant TL software (GE).

Data availability

The data supporting our findings are available in the Supplementary Information, in addition the code is available under (<http://www.bioinf.uni-leipzig.de/publications/supplements/18-059>), RNAseq data are available at the European Nucleotide Archive (ENA; <https://www.ebi.ac.uk/ena>) under the project accession number PRJEB30349, with the following assigned experiment accession numbers.

Notes

1. Here the distance between the 3' end of a gene and the TTS found directly downstream is determined.
2. For better understanding RNA 3' ends identified from the +TEX data are termed TTS throughout the manuscript, all ends identified using the -TEX data are termed RNA 3' ends.

Abbreviations

dRNAseq	differential RNAseq
DSM	Dar-Sorek-Method
IE-PC	Internal Enrichment-Peak Calling
NGS	next-generation sequencing
PS	processing sites
TBP	TATA binding protein
TEX	terminator exonuclease
TFB	transcription factor B
TFE	transcription factor E
TTS	transcription termination site(s)
UTR	untranslated region

Acknowledgments

AM, LKM, JW and PM would like to thank Britta Stoll, Elli Bruckbauer, Irma Merdian and Susanne Schmidt for expert assistance.

Disclosure statement

No potential conflict of interest was reported by the authors.

Funding

This work was supported by the Deutsche Forschungsgemeinschaft [MA1538/21-1] (AM) and by the Austrian Science Fund [SFB F43] (FA).

ORCID

Sarah J. Berkemer  <http://orcid.org/0000-0003-2028-7670>
 Lisa-Katharina Maier  <http://orcid.org/0000-0001-5834-9035>
 Fabian Amman  <http://orcid.org/0000-0002-8646-859X>
 Stephan H. Bernhart  <http://orcid.org/0000-0002-5928-9449>
 Friedhelm Pfeiffer  <http://orcid.org/0000-0003-4691-3246>
 Peter F. Stadler  <http://orcid.org/0000-0002-5016-5191>
 Anita Marchfelder  <http://orcid.org/0000-0002-1382-1794>

References

- [1] Fouqueau T, Blombach F, Cackett G, et al. The cutting edge of archaeal transcription. *Emerging Topics Life Sci.* 2018;2(4):517–533.
- [2] Bell SD. Archaeal transcriptional regulation – variation on a bacterial theme? *Trends Microbiol.* 2005;13:262–265.
- [3] Ray-Soni A, Bellecourt MJ, Landick R. Mechanisms of bacterial transcription termination: all good things must end. *Annu Rev Biochem.* 2016;85:319–347. Epub 062016 Mar 060817.
- [4] Porrua O, Boudvillain M, Libri D. Transcription termination: variations on common themes. *Trends Genet.* 2016;32:508–522. Epub 2016 Jun 1028.
- [5] Peters JM, Vangeloff AD, Landick R. Bacterial transcription terminators: the RNA 3'-end chronicles. *J Mol Biol.* 2011;412:793–813.
- [6] Nemeth A, Perez-Fernandez J, Merkl P, et al. RNA polymerase I termination: where is the end? *Biochim Biophys Acta.* 2013;1829:306–317.
- [7] Arimbasseri AG, Rijal K, Maraia RJ. Transcription termination by the eukaryotic RNA polymerase III. *Biochim Biophys Acta Gene Regul Mech.* 2013;1829:318–330.
- [8] Kuehner JN, Pearson EL, Moore C. Unravelling the means to an end: RNA polymerase II transcription termination. *Nat Rev Mol Cell Biol.* 2011;12:283–294.
- [9] Fouqueau T, Blombach F, Werner F. Evolutionary origins of two-barrel RNA polymerases and site-specific transcription initiation. *Annu Rev Microbiol.* 2017;71:331–348. Epub 092017 Jun 091028.
- [10] Santangelo TJ, Reeve JN. Archaeal RNA polymerase is sensitive to intrinsic termination directed by transcribed and remote sequences. *J Mol Biol.* 2006;355:196–210. Epub 2005 Nov 1019.
- [11] Hirtreiter A, Damsma GE, Cheung AC, et al. Spt4/5 stimulates transcription elongation through the RNA polymerase clamp coiled-coil motif. *Nucleic Acids Res.* 2010;38:4040–4051. Epub 2010 Mar 4042.
- [12] Santangelo TJ, Cubonova L, Skinner KM, et al. Archaeal intrinsic transcription termination in vivo. *J Bacteriol.* 2009;191:7102–7108. Epub 02009 Sep 00911.
- [13] Spitalny P, Thomm M. A polymerase III-like reinitiation mechanism is operating in regulation of histone expression in archaea. *Mol Microbiol.* 2008;67:958–970. Epub 02007 Dec 06019.
- [14] Thomm M, Hausner W, Hethke C. Transcription factors and termination of transcription in methanococcus. *Syst Appl Microbiol.* 1993;16:648–655.
- [15] Santangelo TJ, Cubonova L, Matsumi R, et al. Polarity in archaeal operon transcription in *Thermococcus kodakaraensis*. *J Bacteriol.* 2008;190:2244–2248. Epub 02008 Jan 01811.
- [16] Walker JE, Luyties O, Santangelo TJ. Factor-dependent archaeal transcription termination. *Proc Natl Acad Sci U S A.* 2017;114: E6767–E6773. Epub 1704022017 Jul 1704028131.
- [17] Dar D, Prasse D, Schmitz RA, et al. Widespread formation of alternative 3' UTR isoforms via transcription termination in archaea. *Nat Microbiol.* 2016;1:16143.
- [18] Leigh JA, Albers SV, Atomi H, et al. Model organisms for genetics in the domain Archaea: methanogens, halophiles, Thermococcales and Sulfolobales. *FEMS Microbiol Rev.* 2011;35:577–608.
- [19] Pohlschroder M, Schulze S. *Haloferax volcanii*. *Trends Microbiol.* 2019;27:86–87.
- [20] Ammar R, Torti D, Tsui K, et al. Chromatin is an ancient innovation conserved between Archaea and Eukarya. *eLife.* 2012;1: e00078.
- [21] Babski J, Haas KA, Nather-Schindler D, et al. Genome-wide identification of transcriptional start sites in the haloarchaeon *Haloferax volcanii* based on differential RNA-Seq (dRNA-Seq). *BMC Genomics.* 2016;17:629.
- [22] Bechhofer DH, Deutscher MP. Bacterial ribonucleases and their roles in RNA metabolism. *Crit Rev Biochem Mol Biol.* 2019;54:242–300.
- [23] Pfeiffer F, Broicher A, Gillich T, et al. Genome information management and integrated data analysis with HaloLex. *Arch Microbiol.* 2008;190:281–299. Epub 2008 Jul 2001.
- [24] Veloso F, Riadi G, Aliaga D, et al. Large-scale, multi-genome analysis of alternate open reading frames in bacteria and archaea. *Omic.* 2005;9:91–105.
- [25] Lorenz R, Hofacker IL, Bernhart SH. Folding RNA/DNA hybrid duplexes. *Bioinformatics.* 2012;28:2530–2531.
- [26] Bernhart SH, Hofacker IL, Will S, et al. RNAalifold: improved consensus structure prediction for RNA alignments. *BMC Bioinformatics.* 2008;9:474.
- [27] Lorenz R, Bernhart SH, Zu Siederdisen CH, et al. ViennaRNA Package 2.0. *Algorithms Mol Biol.* 2011;6:26.
- [28] Miladi M, Junge A, Costa F, et al. RNAscClust: clustering RNA sequences using structure conservation and graph based motifs. *Bioinformatics.* 2017;33:2089–2096.
- [29] Afgan E, Baker D, Batut B, et al. The Galaxy platform for accessible, reproducible and collaborative biomedical analyses: 2018 update. *Nucleic Acids Res.* 2018;46:W537–W544.
- [30] Maier LK, Benz J, Fischer S, et al. Deletion of the Sm1 encoding motif in the lsm gene results in distinct changes in the transcriptome and enhanced swarming activity of *Haloferax* cells. *Biochimie.* 2015;6:00058–00059.
- [31] Leon-Sobrinho C, Kot WP, Garrett RA. Transcriptome changes in STSV2-infected *Sulfolobus islandicus* REY15A undergoing continuous CRISPR spacer acquisition. *Mol Microbiol.* 2016;99:719–728.
- [32] Brown JW, Daniels CJ, Reeve JN, et al. Gene structure, organization, and expression in archaeobacteria. *CRC Crit Rev Microbiol.* 1989;16:287–337.

- [33] Koide T, Reiss DJ, Bare JC, *et al.* Prevalence of transcription promoters within archaeal operons and coding sequences. *Mol Syst Biol.* 2009;5:285. Epub 2009 Jun 1016.
- [34] Werner F, Grohmann D. Evolution of multisubunit RNA polymerases in the three domains of life. *Nature Rev Microbiol.* 2011;9:85–98.
- [35] Amman F, D'Halluin A, Antoine R, *et al.* Primary transcriptome analysis reveals importance of IS elements for the shaping of the transcriptional landscape of *Bordetella pertussis*. *RNA Biol.* 2018;15(7):1–9.
- [36] Georg J, Hess WR. cis-antisense RNA, another level of gene regulation in bacteria. *Microbiol Mol Biol Rev.* 2011;75:286–300.
- [37] Thomason MK, Storz G. Bacterial antisense RNAs: how many are there, and what are they doing? *Annu Rev Genet.* 2010;44:167–188.
- [38] Cohen O, Doron S, Wurtzel O, *et al.* Comparative transcriptomics across the prokaryotic tree of life. *Nucleic Acids Res.* 2016;44:W46–53. Epub 2016 May 1096.
- [39] Allers T, Ngo HP, Mevarech M, *et al.* Development of additional selectable markers for the halophilic archaeon *Haloferax volcanii* based on the *leuB* and *trpA* genes. *Appl Environ Microbiol.* 2004;70:943–953.
- [40] Allers T, Barak S, Liddell S, *et al.* Improved strains and plasmid vectors for conditional overexpression of his-tagged proteins in *Haloferax volcanii*. *Appl Environ Microbiol.* 2010;76:1759–1769.
- [41] Martin M. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet J.* 2011;17:10–12.
- [42] Andrews S. FastQC: a quality control tool for high throughput sequence data. Babraham Institute; 2010.
- [43] Hoffmann S, Otto C, Kurtz S, *et al.* Fast mapping of short sequences with mismatches, insertions and deletions using index structures. *PLoS Comput Biol.* 2009;5:e1000502.
- [44] Otto C, Stadler PF, Hoffmann S. Lacking alignments? The next-generation sequencing mapper segemehl revisited. *Bioinformatics.* 2014;30:1837–1843.
- [45] Li H, Handsaker B, Wysoker A, *et al.* The sequence alignment/map format and SAMtools. *Bioinformatics.* 2009;25:2078–2079.
- [46] Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics.* 2010;26:841–842.
- [47] Dar D, Shamir M, Mellin JR, *et al.* Term-seq reveals abundant ribo-regulation of antibiotics resistance in bacteria. *Science.* 2016;352:aad9822.
- [48] Bailey TL, Boden M, Buske FA, *et al.* MEME SUITE: tools for motif discovery and searching. *Nucleic Acids Res.* 2009;37:W202–W208.
- [49] Nawrocki EP, Eddy SR. Infernal 1.1: 100-fold faster RNA homology searches. *Bioinformatics.* 2013;29:2933–2935.
- [50] Holmes ML, Dyall-Smith ML. Sequence and expression of a halobacterial beta-galactosidase gene. *Mol Microbiol.* 2000;36:114–122.
- [51] Delmas S, Shunburne L, Ngo HP, *et al.* Mre11-Rad50 promotes rapid repair of DNA damage in the polyploid archaeon *Haloferax volcanii* by restraining homologous recombination. *PLoS Genet.* 2009;5:e1000552. Epub 1002009 Jul 1000510.
- [52] Large A, Stamme C, Lange C, *et al.* Characterization of a tightly controlled promoter of the halophilic archaeon *Haloferax volcanii* and its use in the analysis of the essential *cct1* gene. *Mol Microbiol.* 2007;66:1092–1106. Epub 2007 Oct 1031.
- [53] Bitan-Banin G, Ortenberg R, Mevarech M. Development of a gene knockout system for the halophilic archaeon *Haloferax volcanii* by use of the *pyrE* gene. *J Bacteriol.* 2003;185:772–778.