

Genome-wide antisense transcription drives mRNA processing in bacteria

Iñigo Lasa^{a,1,2}, Alejandro Toledo-Arana^{a,1}, Alexander Dobin^b, Maite Villanueva^a, Igor Ruiz de los Mozos^a, Marta Vergara-Irigaray^a, Víctor Segura^c, Delphine Fagegaltier^b, José R. Penadés^d, Jaione Valle^a, Cristina Solano^a, and Thomas R. Gingeras^{b,2}

^aLaboratory of Microbial Biofilms, Instituto de Agrobiotecnología, Consejo Superior de Investigaciones Científicas–Universidad Pública de Navarra–Gobierno de Navarra, 31006 Pamplona, Spain; ^bLaboratory of Functional Genomics, Cold Spring Harbor Laboratory, Cold Spring Harbor, NY 11724; ^cGenomics, Proteomics and Bioinformatics Unit, Centro de Investigación Médica Aplicada, Universidad de Navarra, 31008 Pamplona, Spain; and ^dInstituto en Ganadería de Montaña–Consejo Superior de Investigaciones Científicas, 24346 León, Spain

Edited by Susan Gottesman, National Cancer Institute, Bethesda, MD, and approved November 8, 2011 (received for review August 19, 2011)

RNA deep sequencing technologies are revealing unexpected levels of complexity in bacterial transcriptomes with the discovery of abundant noncoding RNAs, antisense RNAs, long 5′ and 3′ untranslated regions, and alternative operon structures. Here, by applying deep RNA sequencing to both the long and short RNA fractions (<50 nucleotides) obtained from the major human pathogen *Staphylococcus aureus*, we have detected a collection of short RNAs that is generated genome-wide through the digestion of overlapping sense/antisense transcripts by RNase III endoribonuclease. At least 75% of sense RNAs from annotated genes are subject to this mechanism of antisense processing. Removal of RNase III activity reduces the amount of short RNAs and is accompanied by the accumulation of discrete antisense transcripts. These results suggest the production of pervasive but hidden antisense transcription used to process sense transcripts by means of creating double-stranded substrates. This process of RNase III-mediated digestion of overlapping transcripts can be observed in several evolutionarily diverse Gram-positive bacteria and is capable of providing a unique genome-wide posttranscriptional mechanism to adjust mRNA levels.

antisense RNA | overlapping transcription | RNA processing | posttranscriptional regulation | microRNA

For many years, the catalog of transcripts (transcriptome) produced by bacterial cells was limited to the transcription products of known annotated genes (mRNA), ribosomal RNA (rRNA), and transfer RNA (tRNA). In the past 10 years, the development of new approaches based on high-resolution tiling arrays and RNA deep sequencing (RNA-seq) has uncovered that a significant proportion (depending on the study, varying between 3% and >50%) of protein coding genes are also transcribed from the reverse complementary strand (1–17). In most cases, overlapping transcription generates a noncoding antisense transcript whose size can vary between various tens of nucleotides (*cis*-encoded small RNAs) to thousands of nucleotides (antisense RNAs). The antisense transcript can cover the 5′ end, 3′ end, middle, entire gene, or even various contiguous genes. Alternatively, overlapping transcription can also be due to the overlap between long 5′ or 3′ UTRs of mRNAs transcribed in the opposite direction. Independent of the mechanism by which it is generated, overlapping transcription has been proposed to affect the expression of the target gene at different levels [for review, see Thomason and Storz (18)]. These mechanisms include: (i) the overlapped transcript affects the stability of the target RNA by either promoting (RNA degradation) or blocking (RNA stabilization) cleavage by endoribonucleases or exoribonucleases; (ii) the overlapped transcript induces a change in the structure of the mRNA that affects transcription termination (transcription attenuation); (iii) the overlapped transcript prevents RNA polymerase from binding or extending the transcript encoded in the opposite strand (transcription interference); and (iv) the overlapped transcript affects protein synthesis either blocking or promoting ribosome binding (translational regulation). Although all these regulatory mechanisms have been proposed based on

studies with specific sense–antisense partners, the presence of massive amounts of overlapping transcription strongly suggest that it might serve for a general purpose on bacterial gene expression (5, 18–24).

In this work, we used RNA sequencing to analyze both the long and short RNA fractions of the major human pathogen *Staphylococcus aureus*. *S. aureus* is a common asymptomatic colonizer of the skin, nasopharynx, and other mucosal surfaces of approximately one-fourth of the healthy human population. However, when *S. aureus* traverses the epithelial barrier, it becomes a leading cause of many diverse pathological syndromes, such as abscesses, bacteremia, endocarditis, osteomyelitis, and pneumonia (25). *S. aureus* has emerged as a model organism for the study of bacterial regulatory RNAs because key discoveries in bacterial regulatory RNAs have been achieved in this bacterium. In 1993, Novick and coworkers (26) identified the first example of a regulatory RNA (RNAIII) that controls the expression of virulence factors by pairing with the target mRNAs followed by degradation of the RNAIII–mRNA complex by the double-stranded specific RNase III (27). More recently, several studies using computational analysis of intergenic regions, microarray technology, and deep sequencing have allowed the identification of >140 small RNAs, including both *trans*- and *cis*-encoded antisense RNAs (10, 28–32). In this current study, we uncover the existence of an overlapping transcription process covering, in a genome-wide extent, the expressed protein coding genes. Base pairing between overlapping RNAs can create double-stranded substrates for RNase III endoribonuclease activity. Such duplex regions promote the cleavage of the double-stranded RNA and the generation of short RNAs (average size of 20 nt). Thus, a collection of stable small RNA molecules that symmetrically map both strands of every region with overlapping transcription is generated. The presence of an identical collection of short RNA molecules that symmetrically mapped both strands of annotated ORFs in *Enterococcus faecalis*, *Listeria monocytogenes*, and *Bacillus subtilis* indicated that this process is evolutionary conserved in Gram-positive bacteria.

Results

Pervasive Antisense Transcription in *S. aureus*. A systematic and hierarchical strategy (Fig. S1) to characterize both long and short

Author contributions: I.L., A.T.-A., and T.R.G. designed research; I.L., A.T.-A., M.V., I.R.d.I.M., M.V.-I., J.R.P., J.V., and C.S. performed research; A.D. and D.F. contributed new reagents/analytic tools; I.L., A.T.-A., A.D., V.S., and T.R.G. analyzed data; I.L., A.T.-A., and T.R.G. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

Data deposition: The sequences reported in this paper have been deposited in the Short Read Archive (accession no. SRP003288.1).

¹I.L. and A.T.-A. contributed equally to this work.

²To whom correspondence may be addressed. E-mail: ilasa@unavarra.es or gingeras@cshl.edu.

This article contains supporting information online at www.pnas.org/lookup/suppl/doi:10.1073/pnas.1113521108/-DCSupplemental.

RNA (<50 nt) fractions from log-phase growing *S. aureus* cells was developed. Long RNA sequencing was performed by using a cDNA synthesis procedure that preserves information about a transcript's direction based on the incorporation of deoxy-UTP during the second strand synthesis and subsequent destruction of the uridine-containing strand (33). The resulting 76-bp paired-end reads were mapped to the *S. aureus* NCTC 8325 reference genome. A total of 9.7 million uniquely mapped read pairs were identified (Fig. S1). 49.2% of the genome was covered by uniquely mapped reads on both strands, 40.4% by uniquely mapped reads on one of the strands, and 10.4% showed no coverage (Fig. 1A). Of the 2,653 annotated ORFs of the *S. aureus* genome, which covers ~84% of the genome, we detected expression of 2,181 ORFs (coverage of >90%), of which 1,387 ORFs displayed 50% coverage on the antisense strand (Fig. 1B).

Naturally occurring short RNAs were also sequenced in a strand-aware fashion by using a two-step adaptor ligation procedure to the 3' and 5' ends of the RNA molecules (34). The reads were aligned by algorithmically clipping off the 3' adaptor, and the remaining sequences of each read were mapped to the genome by using STAR (<http://gingeraslab.cshl.edu/STAR/>). For

alignments of 10–19 bases long, up to one mismatch was allowed; for alignments >20 bases, up to two mismatches were allowed. Alignments of <10 bases were discarded, and spliced alignments were prohibited. This process yielded a total of 7,778,726 million reads mapped to the genome (Fig. S1). The average length of short RNA molecules was 20 nt (Fig. S1). The uniquely mapped short RNA sequences covered, in at least 50% of their length, 2,268 and 1,981 ORF regions on the sense and antisense strands, respectively (Fig. 1C). Thus, the percentage of ORFs covered in at least 50% of their length by reads in the antisense strand was higher in the case of short RNA (75%) than in the case of long RNA (56%), suggesting that short RNA libraries may prove to be a more sensitive way to detect antisense transcripts. Overall, and given that long and short RNA libraries were generated independently—that is to say from two fractions coming from the same RNA sample—these results provide evidence of the existence of antisense transcription not seen in the long RNA sequence analysis.

Symmetric Distribution of Short RNA Reads in both Strands of ORF Regions.

We next sought to determine whether the distribution of short and long reads for a given ORF were linked. For that, we visualized normalized \log_2 values representing the number of mapped reads per nucleotide using the Integrated Genome Browser (IGB) (35). Fig. 2 illustrates a randomly selected 30-kb region of the genome of *S. aureus* that represents 1% of the genome and depicts the uniquely mapped long and short RNAs. The results revealed that short RNA sequences were symmetrically distributed in both strands of the ORFs, whereas long RNA transcripts follow the expected biased distribution toward the sense strand. Intriguingly, the regions with detectable overlapped transcription between long RNA transcripts, such as those regions corresponding to antisense transcripts to ORFs (00056, 00061, *sirABC* operon), were covered with higher numbers of short RNA reads in both strands. Similar symmetrical accumulation of high levels of short RNAs was detected in every region of the genome where noticeable overlapping transcription occurs, such as 5' and 3' overlapping UTRs, overlapping operons (ORFs that, being located in the middle of an operon, are transcribed in opposite direction to the other genes of the operon), and antisense transcripts (see Fig. 3 and Figs. S2–S4 for additional examples). To most accurately demonstrate that the distribution of short RNA reads was symmetric genome-wide, we quantified the number of long and short RNAs mapping to the sense and antisense strands of each ORF. In accordance with the images observed with the IGB browser, the results revealed very similar numbers of short RNA reads genome-wide in both strands of ORF regions and the expected biased distribution of long RNA reads in the sense strand (Fig. 4A and B). In summary, these results show that the *S. aureus* transcriptome contains both long and very short RNA molecules. The amount of long RNAs is, as expected, higher in the sense strand of each ORF. In contrast, short RNAs are equally distributed in both strands of each ORF and specially enriched in those regions with detectable overlapped transcription between long RNAs.

RNase III Is Responsible for the Production of Symmetrically Distributed Short RNA Populations.

The fact that short RNAs display a symmetrical distribution in sense/antisense strands and accumulate in higher numbers in regions with noticeable overlapping transcription raised the possibility that short RNA molecules were derived from the cleavage of overlapping long sense/antisense primary transcripts. *S. aureus* genome has been reported to encode at least eight putative endoribonucleases and three exoribonucleases (32). Among them, the RNase III endoribonuclease is the only enzyme known to be able to degrade double-stranded RNA. Thus, we tested the possibility that RNase III might be responsible for processing the overlapping transcripts into symmetrically distributed sense and antisense short RNA populations. An RNase III mutant in the *S. aureus* 15981 background (*S. aureus* 15981 Δ rnase III) was constructed by using a described approach (36).

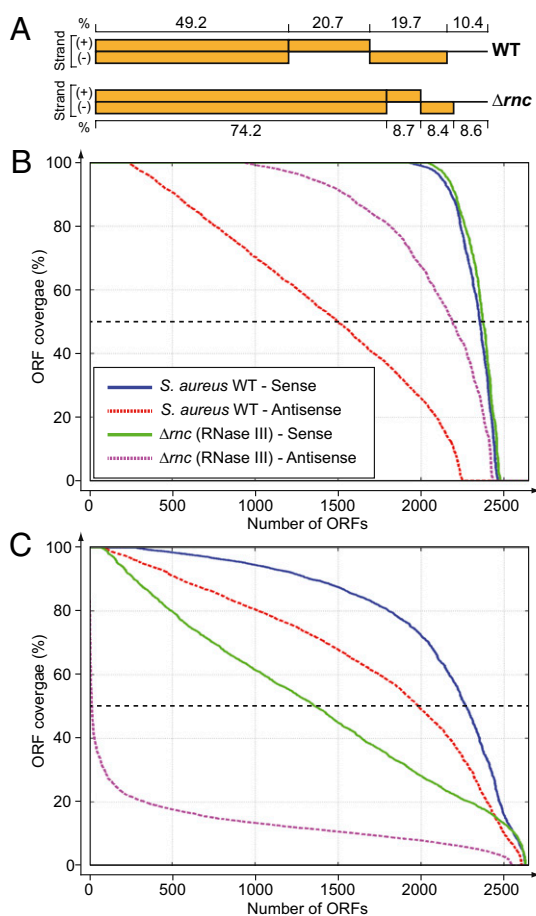


Fig. 1. Genome-wide analysis of mapped reads from long and short RNA-seq libraries. (A) Percentage of the genome of *S. aureus* NCTC 8325 covered by uniquely mapped reads on both strands, reads on one of the strands, and showed no coverage, respectively. The long RNA-seq libraries were prepared from *S. aureus* 15981 wild-type strain (WT) and its corresponding RNase III mutant (Δ rnase III). (B and C) Comparison of the cumulative distribution of ORF coverage by long (B) and short (C) RNA reads. The plot represents the number of ORFs (x axis) found above the ORF coverage value (y axis). The coverage was computed from the collapsed reads uniquely mapped in the sense and antisense orientation to the ORFs. The dashed line represents 50% coverage.

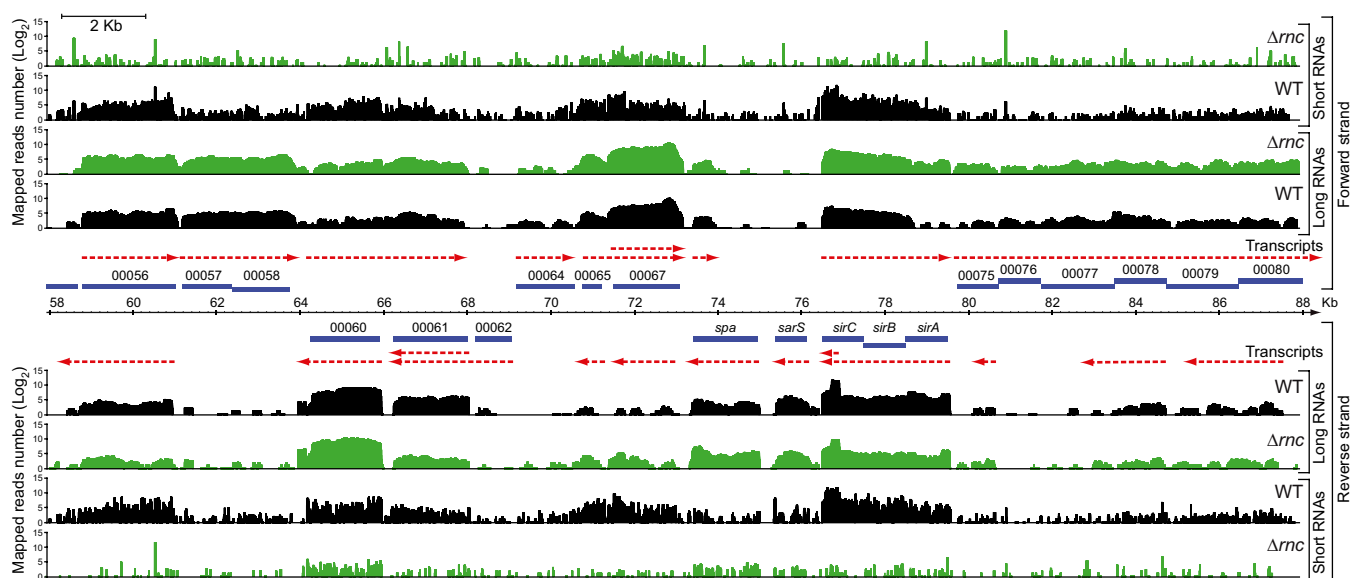


Fig. 2. Long and short mapped reads distribution in *S. aureus* genome. The drawing is an IGB software image showing the uniquely mapped long and short RNAs in a 30-kb region (1%) of the genome of *S. aureus* NCTC 8325. Transcripts are represented as dashed red arrows. Genomic coordinates denote the position in kilobases of the *S. aureus* NCTC 8325 genome. Annotated ORFs are shown as blue lines. The number on the ORF indicates the gene identification. Long and short RNAs show the distribution of uniquely mapped reads of long and short RNA libraries. *S. aureus* 15981 (black) and *S. aureus* 15981 Δrnc (RNase III mutant) (green). The scale (\log_2) indicates the number of mapped reads per nucleotide position.

Analysis of the uniquely mapped reads from long RNA-seq of the RNase III mutant revealed that the percentage of the genome covered by reads on both strands increased up to 74.2% compared with wild-type strain (49.2%; Fig. 1A). This increase was mainly due to a significantly higher coverage of the antisense strand (82% ORFs displayed 50% coverage on the antisense strand; Fig. 1B). In contrast, the number of the short RNA reads was drastically reduced, especially in the antisense strand (Fig. 1C), reducing the percentage of the genome that was covered on both strands by short RNAs to only 6%. Of note, the median length of the short RNA molecules in RNase III mutant was 15 nt, suggesting the possibility that short RNAs detected in the RNase III mutant were produced by another RNA processing pathway (Fig. S1). Visualization of the distribution of mapped reads by using IGB confirmed that short RNA sequences had lost their symmetric distribution, whereas long RNA transcripts followed the expected biased distribution toward the sense strand (Fig. 2 and Figs. S2–S4). Accordingly, the correlation first observed in the wild-type strain between the numbers of short RNA reads in sense and antisense strands per annotated ORF region disappeared in the analysis of the RNase III mutant (Fig. 4C and D). Together, these results indicate that a large majority of short RNA molecules present in the transcriptome of *S. aureus* are produced by the cleavage activity of double-stranded RNase III enzyme.

Because the pattern and cleavage frequency by RNase III is unknown, the short RNA molecules might be direct products of RNase III activity or processed products of larger RNA fragments generated by RNase III. Pnp is the most important 3'–5' exonuclease activity in bacteria. *S. aureus* contains a gene (SAOUHSC_01251) encoding a protein that shares 66% identity with Pnp of *B. subtilis*. We produced libraries from short RNA fraction of *S. aureus* 15981 Δpnp . Analysis of the mapped reads from *S. aureus* 15981 Δpnp mutant revealed that the distribution and size of the short RNAs followed a pattern indistinguishable from that of the wild-type strain (Fig. S1), suggesting that Pnp activity is not required for subsequent processing of the short RNA molecules generated by RNase III activity.

Abundance of Short RNAs Correlates with the Levels of Double-Stranded Sense/Antisense Transcripts. One prediction of the model that short RNAs are produced from the processing of genome-

wide overlapping regions of transcription is that the abundance of short RNAs detected should be proportional to the abundance of either the sense/antisense transcripts, depending upon which transcribed strand is less abundant and available for processing. To explore this prediction, we analyzed the short and long RNA complements from a sigma B (*sigB*) mutant (*S. aureus* $\Delta sigB$). The transcription factor Sigma B drives the transcription activity of genes under specific environmental conditions. We analyzed ORF regions for which the long antisense transcripts contained a consensus SigB promoter box (Fig. S5), and their expression was significantly suppressed in the *sigB* mutant (>50% reduction in the $\Delta sigB$ /WT antisense transcript ratio based on the long RNA libraries). Consistent with the hypothesis that the abundance of short RNAs depends on the levels of double-stranded RNA, knockdown of the antisense transcripts levels in *sigB* mutant correlates with a decrease in the amount of short RNAs produced from both strands (Fig. S5). These results indicated that the short RNA abundances at ORF regions are strongly correlated with the less abundant levels of overlapping long RNA capable of forming double-stranded RNA.

Detection of Occurrence and Abundance of Antisense Transcripts in RNase III Mutant. Detection of antisense transcripts has been difficult in bacteria, and only the presence of a few antisense transcripts has been confirmed by Northern blot techniques. RNase III cleavage of overlapping RNA transcripts into short RNA molecules could explain at least in part the paucity of antisense transcripts detected so far. Only minimal amounts of antisense transcripts would be maintained in the cell by RNase III activity. To explore this hypothesis, we performed Northern hybridizations with strand-specific probes to interrogate sense and antisense transcripts of several individual genes in wild-type and RNase III mutant strains. The candidate genes were selected based on their relevance to different aspects of *S. aureus* virulence (*sarA*, *agrBCDA*, *saePQRS*, *clpP*) or biology (*lexA*, *recF*, *yhSR*) (Fig. S6). The results of the Northern analyses indicated a specific absence and presumed degradation of most full-length antisense transcripts in the steady-state condition of the wild-type strain (Fig. 5), whereas the presence of discrete-size antisense transcripts was clearly detectable in the RNase III mutant for all genes tested. It is worth noting that these results confirm

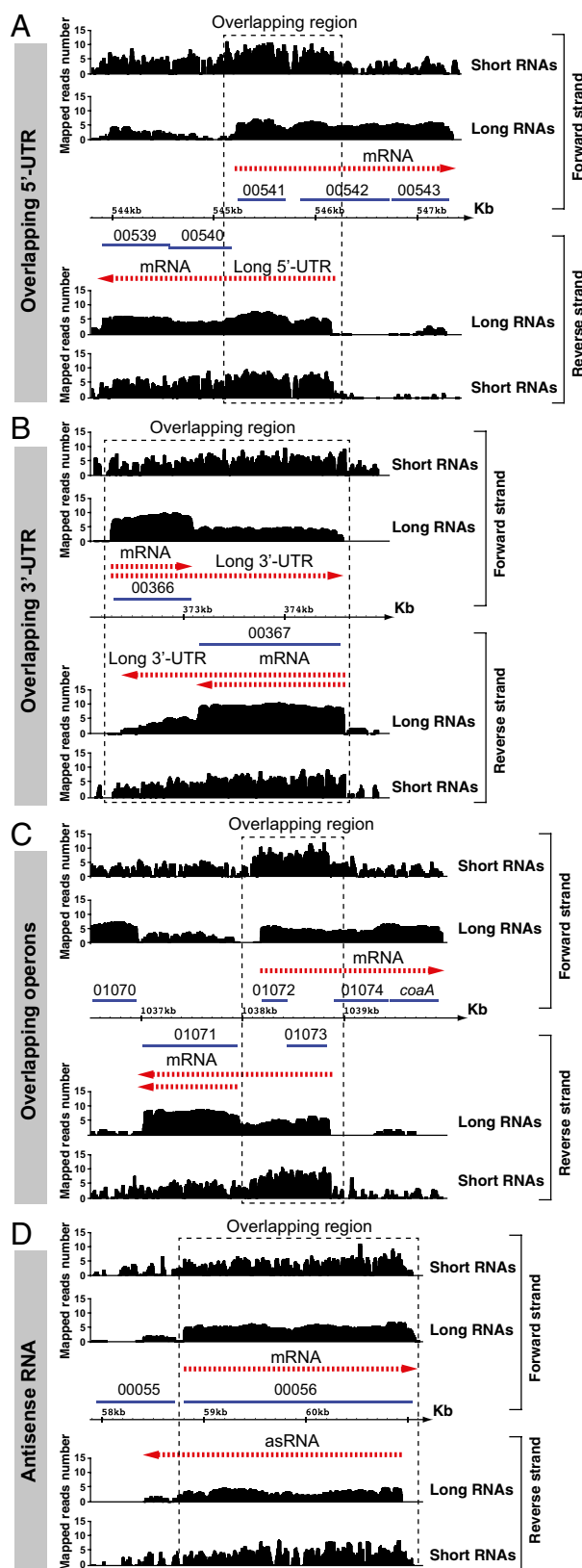


Fig. 3. Examples of mapped reads distribution in regions with overlapping transcription of *S. aureus*. Drawings are images from IGB software showing different regions of the genome of *S. aureus* NCTC 8325. Examples of overlapping 5' UTRs (A), overlapping 3' UTRs (B), overlapping operons (C), and antisense RNA (D) are shown. Transcripts are represented as dashed red arrows. Genomic coordinates denote the position in kilobases of the *S. aureus*

the existence of antisense transcripts for genes that have been thoroughly studied because of their impact on *S. aureus* virulence and antibiotic resistance. For some genes, these hybridizations showed that the RNA levels of the sense strand (*lexA*, *clpP*, *saePQRS*) increased in the RNase III mutant, suggesting that the absence of RNase III cleavage can slightly modulate the expression levels of sense transcripts (Fig. 5). To confirm that the presence of antisense RNA was restricted to those regions where short RNAs were detected, we selected two genes (SAOUHSC_00086 and SAOUHSC_00410) for which very few short RNAs were detectable in the wild-type strain (Fig. S7). In both cases, no transcript antisense to these genes was detectable in the RNase III mutant. Overall, these results uncover the existence of long antisense RNAs transcripts for most ORFs of the *S. aureus* genome. These long antisense transcripts are underrepresented in the wild-type strain because of the double-stranded RNase activity of RNase III.

Analysis of Short RNA Complement Present in Diverse Bacterial Species. To investigate whether this genome-wide sense/antisense overlapping transcript processing mechanism is specific to *S. aureus* or is active in other bacterial species, we characterized the short RNA complement present in three representative Gram-positive (*E. faecalis*, *L. monocytogenes*, and *B. subtilis*) and one Gram-negative (*Salmonella enterica* serovar Enteritidis) bacteria (Fig. 4E). Short RNA libraries were produced, sequenced, and mapped by using the described protocol (Fig. S1). Analysis of the distribution of short RNAs in sense and antisense strands of ORF regions revealed a highly significant correlation between quantities of short RNAs in sense/antisense strands for the three low-GC content Gram-positive bacteria, mirroring the observations in *S. aureus*. In contrast, the results obtained from the analysis of *Salmonella* demonstrated the absence of such a correlation, indicating the existence of a different processing pattern of overlapping RNA pairs in Gram-negative bacteria. Previous transcriptome analysis has allowed the identification of antisense transcripts in *L. monocytogenes* (5), *B. subtilis* (17), and *E. faecalis* (37). Analysis of the distribution of short RNAs in those regions with recognized antisense transcription confirmed the accumulation of high amounts of short RNA in every region with overlapping transcription, indicating that genome-wide digestion of overlapping sense/antisense transcripts is conserved at least in Gram-positive bacteria (Fig. S8).

Discussion

Development of RNA-seq technology is allowing the characterization of the multiple types of RNA molecules present in a living cell. The application of this technology in bacteria has primarily been restricted to the analysis of long RNA molecules because of the difficulty in removing highly abundant small-sized rRNA and tRNA molecules. Here, we have used two methods developed for microRNA analysis in eukaryotic cells to analyze the RNA fraction of <50 nt of the human pathogen *S. aureus*. The short RNA fraction was purified by size fractionation electrophoresis, and libraries for RNA-seq were generated by following a protocol that preserves the information about a transcript's direction developed for the direct cloning of microRNA in *Drosophila* (34).

The analysis of the distribution of short RNA molecules revealed several unexpected results. First, the sense strand of 2,268 ORFs and the antisense strand of 1,981 ORFs were covered with unique short RNA reads in at least 50% of their length, indicating the existence of antisense transcription from

NCTC 8325 genome. Annotated ORFs are shown as blue lines. The number on the ORF indicates the gene identification. Long and short RNAs show the distribution of uniquely mapped reads of long and short RNA libraries in *S. aureus* 15981. The scale (\log_2) indicates the number of mapped reads per nucleotide position. Dashed rectangles highlight increased accumulation of short mapped reads in regions with overlapping transcription, according to long RNA reads.

