

Dynamics of transcriptional start site selection during nitrogen stress-induced cell differentiation in *Anabaena* sp. PCC7120

Jan Mitschke^a, Agustín Vioque^b, Fabian Haas^a, Wolfgang R. Hess^{a,1} and Alicia M. Muro-Pastor^b

^aUniversity of Freiburg, Faculty of Biology, Schänzlestr. 1, D-79104 Freiburg, Germany

^bInstituto de Bioquímica Vegetal y Fotosíntesis, Consejo Superior de Investigaciones Científicas and Universidad de Sevilla, E-41092 Seville, Spain

¹Corresponding author: Wolfgang R. Hess; Phone: +49-761-203-2796

Fax: +49-761-203-2445; Email: Wolfgang.Hess@biologie.uni-freiburg.de

Keywords: prokaryotic cell differentiation, antisense transcription, heterocyst, nitrogen fixation, differential RNA-seq

Abbreviations: asRNA, antisense RNA; dRNA-seq, differential RNA-seq; HEP, heterocyst envelope polysaccharides; ncRNA, noncoding RNA; PSSM, position-specific scoring matrix; TAP, tobacco acid pyrophosphatase; TEX, TerminatorTM 5'phosphate-dependent exonuclease; TSS, transcriptional start site

1 ABSTRACT

2 The fixation of atmospheric N₂ by photosynthetic cyanobacteria is a major source of
3 nitrogen in the biosphere. In Nostocales, such as *Anabaena* sp. PCC7120, this
4 process is spatially separated from oxygenic photosynthesis and occurs exclusively
5 in heterocysts. Upon nitrogen step-down, these specialized cells differentiate from
6 vegetative cells in a synchronized process that is controlled by two major regulators:
7 NtcA and HetR. However, the regulon controlled by these two factors is only partially
8 defined, and several aspects of the differentiation process have remained enigmatic.
9 Using dRNA-seq, we experimentally define a genome-wide map of more than 10,000
10 transcriptional start sites (TSS) of *Anabaena* sp. PCC7120, a model organism for the
11 study of prokaryotic cell differentiation and N₂ fixation. By analyzing the adaptation to
12 nitrogen stress, our global TSS map provides insight into the dynamic changes that
13 modify the transcriptional organization at a critical step of the differentiation process.
14 We identify more than 900 TSS with an absolute fold change in response to nitrogen
15 deficiency of at least eight. From these, at least 209 were under control of HetR,
16 whereas at least 158 other TSS were not, but potentially directly controlled by NtcA.
17 Most of the promoters activated during the switch to N₂ fixation were previously
18 unknown, thereby adding hundreds of protein-coding genes and non-coding
19 transcripts to the list of potentially involved factors. These data experimentally define
20 the NtcA regulon and the DIF+ motif, a palindrome at or close to position -35 that
21 seems essential for heterocyst-specific activation of expression of certain genes.

22

1 \body

2 **Introduction**

3 Cyanobacteria are oxygen-producing, photosynthetic organisms that are responsible
4 for approximately half of the global CO₂ fixation. In addition, many cyanobacteria are
5 able to perform N₂ fixation, a process that is extremely sensitive to oxygen. To
6 protect the nitrogenase complex from photosynthetically-evolved oxygen, many
7 filamentous strains, including *Anabaena* sp. PCC7120 (also known as *Nostoc* sp.
8 PCC7120, from here on *Anabaena* 7120), differentiate heterocysts, a specialized cell
9 type devoted to N₂ fixation (1, 2). Heterocyst differentiation is integrated into a series
10 of physiological responses that take place when a source of combined nitrogen is not
11 available. Those responses are globally controlled by NtcA, a transcriptional
12 regulator of the cyclic AMP receptor protein (CAP) family (3-5). NtcA-mediated
13 regulation involves the binding of NtcA to a consensus binding site with the sequence
14 GTAN₈TAC (3). In the absence of ammonium, the preferred nitrogen (N) source,
15 NtcA (5) activates the expression of genes required for alternative assimilation
16 pathways, such as those encoding nitrate and nitrite reductases (3). NtcA also acts
17 as a transcriptional repressor of some genes, such as the *gif* gene, which encodes
18 the glutamine synthetase inactivating factor (5). NtcA is also required for heterocyst
19 differentiation and subsequent N₂ fixation (6, 7). A key factor activated by NtcA is
20 HetR, a master regulator of many genes involved in the differentiation process (8, 9).
21 HetR binds to DNA (10) and folds into an unusual structure (11). Although HetR was
22 hypothesized to control the expression of hundreds of genes, only a single 17-base
23 pair palindrome has been identified as a binding site (12). Moreover, HetR exerts
24 positive feedback on NtcA expression (13), but it is not known how the double-
25 positive feedback between both factors is terminated at a later step in the
26 differentiation cascade.

1 During the last three decades, fewer than 100 transcriptional start sites (TSS)
2 have been mapped in *Anabaena* 7120 on a gene-by-gene basis, primarily associated
3 with highly expressed genes or genes regulated by nitrogen availability. An analysis
4 of the *Anabaena* 7120 transcriptome during vegetative cell growth and in response to
5 N deprivation has been recently published (14). However, the RNA samples were not
6 enriched for primary transcripts, causing ~90% of the reads at each time point to
7 correspond to processed rRNA and tRNA transcripts and rendering the direct
8 identification of TSS difficult (14).

9 To complement the scarce existing information, our study utilized a differential
10 RNA-seq (dRNA-seq) approach, which is selective for the 5' ends of primary
11 transcripts (15, 16) and allows the comprehensive determination of the transcriptional
12 organization of a genome. Pretreatment of bacterial RNA with TerminatorTM 5'
13 phosphate-dependent exonuclease (TEX) specifically degraded transcripts with a 5'
14 P (processed RNAs) (16). Based on this approach, we present a genome-wide map
15 of 13,705 candidate TSS that were experimentally mapped for the chromosome and
16 the six plasmids of *Anabaena* 7120. Analyzing the transcriptional changes that occur
17 upon nitrogen step-down in both the wild-type (WT) and a *hetR* mutant lead to the
18 identification of all promoters controlled directly or indirectly by HetR, to a precise
19 definition of the NtcA regulon and to the DIF+ motif, a previously unnoticed sequence
20 element involved in heterocyst-specific expression. The availability of its annotated
21 primary transcriptome will greatly facilitate the use of this genetically tractable
22 organism as a model for prokaryotic cell differentiation and N₂ fixation.

23

24 **Results**

25 **Large-scale mapping of primary 5' ends using dRNA-seq**

1 RNA samples were obtained from WT and *hetR* mutant strains that were grown in
2 the presence of ammonia (WT-0 and *hetR*-0) or subjected to N step-down for 8 h
3 (WT-8 and *hetR*-8). In total, 24,312,062 sequence reads (up to 75 nt long) were
4 analyzed, and 1.6 billion bases of cDNA were mapped to the *Anabaena* 7120
5 chromosome and its six large plasmids. In the four samples, between 48.3 and
6 64.7% of the reads did not correspond to rRNAs. The complete dataset was taken to
7 identify possible TSS, based on a minimum number of 50 sequencing reads
8 associated with an RNA 5' end. An example of the data obtained is shown for the
9 gene encoding HetR in **Fig. 1A**. All four previously described TSS, including that at -
10 271 whose induction is heterocyst-specific (13, 17, 18), were identified by a total of
11 6,900 reads in our dataset.

12 We identified 12,797 putative chromosomal TSS and 908 putative TSS on the
13 six plasmids (**Table 1**). From these, 4,186 TSS were located within a distance of 200
14 nt upstream of an annotated gene (gTSS, mostly mRNAs); 4,172 TSS in inverse
15 orientation (or ≤ 50 bp 5' or 3') to annotated genes (aTSS), suggesting antisense
16 transcription; and 1,414 TSS of potential ncRNAs in intergenic spacers (nTSS). In
17 addition, 3,933 TSS in sense orientation were located internally within annotated
18 genes (iTSS). For consistency, this classification (**Fig. S1A**), solely based on
19 location, was used throughout. Therefore, some of the TSS here categorized as
20 gTSS may actually give rise to ncRNAs and some of the nTSS may rather drive the
21 transcription of genes with long 5'UTRs (see also comments in **Table S1**). Because
22 for 704 TSS an association with more than one category was possible (**Fig. S1B**), we
23 prioritized gTSS over aTSS and iTSS and all remaining TSS were automatically
24 categorized as nTSS. A global overview of the distribution of TSS is given in **Fig. 1B**
25 (chromosome) and **Fig. S2** (plasmids). The exact positions of all putative TSS are
26 indicated in **Supplementary Data Files 1** (chromosome) and **2-7** (plasmids) and in

1 **Table S1.** To benchmark, we compared our data with a set of 93 TSS previously
2 reported for 60 different genes or operons from 59 independent studies (**Table S2**).
3 From the previously reported TSS, 81 were confirmed with 69 of these being
4 associated with 50 or more reads. In addition, TSS that had remained unnoticed
5 were observed for several genes, including *glnA*, *ntcA*, *patS* and *rbcL*.
6 Finally, although processing of tRNA precursors is fast, 37 TSS were identified
7 unambiguously for 35 tRNA genes. Based on their alignment, a consensus for a
8 constitutive promoter was defined (**Fig. S3**). We observed that the length of the 5'
9 leaders varied from 5 to >200 nt, but most were between 10 and 20 nt (**Fig. S3 and**
10 **Supplementary data file 8**).

11

12 **Nitrogen deficiency- versus differentiation-related promoters in the nitrogen** 13 **stress response of *Anabaena* 7120**

14 To analyze the transcriptional changes induced by nitrogen stress and those
15 specifically leading to the differentiation of heterocysts, we individually compared the
16 numbers of dRNA-seq reads of all TSS identified in RNA isolated from the four
17 different samples. Normalized ratios (fold changes) in the number of reads between
18 the different samples were determined according to established protocols (19).

19 The comparison of previously described N stress-induced transcriptional
20 responses, e.g., for the *nirA-nrtABCD-narB* cluster (NtcA-activated), the *gifA* gene
21 (NtcA-repressed), and genes in the heterocyst envelope polysaccharide (HEP)
22 island, which are involved in heterocyst maturation, revealed a high degree of
23 consistency (**Fig. S4**). We confirmed the TSS for *nirA* (-460 with respect to the
24 translational start) (20), with a number of reads significantly higher in the samples
25 collected 8 h after N step-down from both the WT and the *hetR* mutant (**Fig. S4A**).
26 We also confirmed the TSS for *gifA* at position -43 (21), associated with a high

1 number of reads in the presence of NH_4^+ from both the WT and the *hetR* mutant (**Fig.**
2 **S4B**). **Fig S4C** shows that transcription of several genes involved in the synthesis of
3 HEP was almost completely restricted to WT under N stress.

4 Two main groups of TSS were defined that differed in their response to N
5 step-down. The DEF category (deficiency-related changes) includes TSS showing
6 transcriptional changes common to both strains (TSS for *nirA* or *gifA* above are
7 paradigms for this category), whereas the DIF category (differentiation-related
8 changes) includes TSS with transcriptional changes observed exclusively in the WT
9 (e.g., TSS for genes in the HEP island, **Fig S4C**). The DIF group includes all
10 transcriptional changes that depend on HetR and are thus likely involved in the
11 process of heterocyst differentiation. Additionally, although most changes in the DEF
12 category involved activation (DEF+), we also identified some TSS associated with a
13 decrease in the number of reads upon N step-down (DEF-) (**Tables S3 to S5**). With a
14 minimum fold change of eight, we identified 129, 28 and 209 TSS in the DEF+, DEF-
15 and DIF+ categories, respectively. Our dataset identifies, for the first time, strongly
16 regulated TSS for many genes with previously described N-dependent regulation or
17 role in heterocyst differentiation, including *amt1* and *amt4*, heterocyst differentiation-
18 related genes *hepA*, *hepB*, *hepN*, and *hepS*, regulatory genes *nirB* and *patB*, or *nbIA*
19 involved in phycobilisome degradation (**Table S6**).

20 Our dataset allows the identification of multiple TSS in complex promoters
21 (**Fig. 1A**), providing a powerful approach to the analysis of genes with complex
22 regulation. Two such promoter regions containing several previously unidentified
23 TSS were chosen for further validation by primer extension or northern blot
24 hybridization. Five putative TSS were identified for *nbIA* (**Table S6; Fig. 2A**), an N
25 stress-inducible gene required for phycobilisome degradation but not essential for
26 heterocyst differentiation in *Anabaena* 7120 (22). The results presented in **Fig. 2B**

1 confirmed all 5' ends identified by dRNA-seq and their activity: TSS1 was not active
2 in the presence of ammonium and NtcA-dependent, TSS2 and TSS5 were also
3 inducible but NtcA-independent, TSS3 was barely detected, and TSS4 was mostly
4 detected in the *ntcA* mutant. Two TSS were identified for *alr3808* (**Table S6; Fig. 2C**)
5 encoding a DpsA homologue with known N-dependent regulation (23, 24).
6 Consistent with the dRNA-seq data, two transcripts covering *alr3808* were identified
7 by northern blot (**Fig. 2D**). The longer transcript, probably originating at position
8 4601709f, is induced upon N step-down but was not expressed in the *hetR* mutant
9 (therefore categorized as DIF+), whereas the shorter transcript, probably originating
10 at position 4601982f, was also induced in the *hetR* mutant, although at later time
11 points. Thus, not only the positions of TSS but also the regulation observed by
12 dRNA-seq were confirmed by different methodologies.

13

14 **TSS in the DIF category identify novel HetR-regulated elements and a sequence** 15 **motif related to heterocyst-specific expression**

16 HetR is the earliest known dedicated regulator involved in the differentiation of
17 functional heterocysts. Therefore, transcriptional responses observed in the WT but
18 not in the *hetR* mutant probably belong to the specific transcriptional program that
19 leads to the differentiation of these cells. These TSS constitute the DIF+ category
20 (see TSS with at least 8-fold change in **Table S5**). Their promoters can be analyzed
21 to identify elements that might play a role in the differentiation process. In fact, the
22 TSS for some genes with known HetR-dependent or heterocyst-specific expression
23 (e.g. *ntcA*, *hetR*, *nsiR1*) appear among the TSS exhibiting the highest fold change in
24 this category. A direct search for the 17-base pair palindrome identified as a HetR
25 binding site (12), or for any other conserved element was unsuccessful. However, we
26 noticed that the promoters for ncRNA NsiR1 (25) were in the DIF+ class. Because

1 NsiR1 is conserved in Nostocales and transcribed from a tandem array of short
2 repeats, we could compare the promoter regions of 51 repeats from five different
3 strains and found a conserved palindrome 5'TCCGGA at or close to the -35 position.
4 Moreover, the same or a very similar motif is present in several other heterocyst-
5 specific promoters (**Fig. 3A**). A global search identified this motif at similar position in
6 58 of the 209 DIF+ promoters (**Table S5**; selected examples in **Fig. 3A**) when a
7 single mismatch was allowed (**Table S7**). From all remaining 13,496 TSS only 572
8 also share this motif (and some of those are DIF+ too, but with a fold change <8).
9 Hence the enrichment for this motif within the DIF+ category of promoters is non-
10 random ($P < 2.2e^{-16}$ in a Chi-squared test). We therefore named this sequence the
11 DIF+ motif.

12 To directly test the functional relevance of the DIF+ motif, the 70 bp promoter
13 from NsiR1 repeat 6 (P6) was placed upstream of a promoter-less GFP. Constructs
14 bearing the intact DIF+ motif expressed green fluorescence exclusively in
15 (pro)heterocysts (**Fig. 3B**). Hence P6 possesses all the elements for cell-specific
16 expression. In contrast, only very weak, non-heterocyst-specific fluorescence was
17 obtained when the DIF+ motif was replaced by 5'GAATTC (**Fig. 3B**). Thus the DIF+
18 motif is required for heterocyst-specific expression of the *nsiR1* promoter.

19

20 **The NtcA binding site revisited: TSS in the DEF category define the NtcA** 21 **regulon**

22 We hypothesized that most N step-down-induced responses occurring both in the
23 *hetR* strain and the WT (DEF categories) would not be related to heterocyst
24 differentiation but rather are likely to be NtcA-regulated. The NtcA binding sites in the
25 NtcA-activated promoters overlap in most cases the -35 region and are centered
26 close to position -41.5 (i.e., the first nt is located at -48) with regard to the TSS (3).

1 **Fig. 4A** shows the promoters of the 20 TSS exhibiting the highest fold change in the
2 DEF+ category. 18 of them contain sequences matching the consensus NtcA binding
3 site at the expected position. Whereas NtcA-dependent activation was previously
4 described for two of them, 580293f (26) and 5167792r (27), the remaining 16 strongly
5 regulated TSS are novel. Thus, the DEF category defines the NtcA regulon at an
6 unprecedented resolution. A position-specific scoring matrix (PSSM) was defined
7 based on an alignment of all promoter regions for TSS in the DEF+ category with at
8 least 8-fold change (**Fig. S5A**). In addition to the expected conservation of positions
9 1-3 and 12-14, this analysis indicated that G residues are strongly avoided at
10 positions 5, 7 and 11 of the NtcA binding site, whereas A/C or A/T are somewhat
11 preferred at positions 5-7 (**Fig. 4B, C**). When all TSS in the DEF+ category were
12 scanned in a sliding window approach for possible NtcA binding sites with a score
13 ≥ 5 , a peak was identified at positions -48 to -49 (first nt of the motif), thereby
14 matching the expected location precisely not only for this dataset with ≥ 8 -fold change
15 (**Fig. 4D**), but also when a much larger dataset of 965 TSS with ≥ 2 -fold change (**Fig.**
16 **S5**) was analyzed. This observation indicated a significant proportion of the DEF+
17 TSS (even with relatively low fold changes) indeed contain NtcA binding sites at
18 positions that are compatible with transcriptional activation and that the PSSM as
19 defined here (**Fig. 4 and S5**) is useful for global searches.

20 NtcA binding sites incompatible with transcriptional activation (eventually
21 repressing transcription) can be located closer to, and even downstream of, the TSS.
22 To find such elements, we scanned sequences surrounding all TSS in the DEF
23 category in two windows: position -120 to -44 (first nt of motif; activation-compatible
24 sites), and -44 to +41 (repression-compatible sites) and show sites with score ≥ 5 in
25 **Tables S8 and S9. Fig. 4E** shows ten examples for putative NtcA binding sites at
26 repression-compatible positions around TSS in the DEF- category. Two of these

1 TSS, 2809313r (21) and 2807328f (28), were previously described as NtcA-
2 repressed. TSS 1785466f (*rbcL*), also described as containing an NtcA binding site in
3 a repressor-compatible position (29), is included for comparison.

4

5 **Novel non-coding RNAs potentially involved in the response to N stress and** 6 **heterocyst differentiation**

7 Several of the TSS exhibiting the highest fold change in our dataset correspond to
8 antisense (asRNA) or non-coding (ncRNA) transcripts. We have further confirmed
9 transcription from some of them, including their regulation. The aTSS at position
10 3953418f, a strongly regulated DEF+ promoter, gives rise to an asRNA for gene
11 *a//3278*, whose mutation leads to the inability to fix N₂ in the presence of oxygen (30).
12 Primer extension analysis confirmed the dRNA-seq results (**Fig. 5A**). The initiation of
13 transcription at this position is strongly induced by N step-down, independently of
14 HetR, but depending on NtcA, consistent with the identification of a putative NtcA
15 binding site upstream (**Table S8; Fig. 4A**).

16 We also confirmed two strongly regulated nTSS that produce small ncRNAs
17 (**Fig. 5B**). Transcription from position 3141905r (DIF+) produces NsiR2, whereas
18 transcription from position 5452083f (DEF+) produces NsiR3. The co-regulation of
19 these TSS with well-studied protein-coding genes in these categories suggests that
20 some of the ncRNAs identified here might be involved in the adaptation to N-stress or
21 the differentiation of functional heterocysts.

22

23 **DISCUSSION**

24 In this study, we have defined a set of more than 10,000 putative TSS for *Anabaena*
25 7120. We did not require these 5' ends to be linked to a classical -10 element
26 because the differentiation process could involve alternative sigma factors

1 recognizing different promoter elements, because the quality of data appeared high
2 (3,401 TSS were identified on the basis of more than 300 reads) and because our
3 dataset was experimentally validated in several ways. This dataset confirms most of
4 the previously defined TSS for this organism (**Table S2**), while identifying new N-
5 regulated TSS for the majority of genes previously reported as involved in heterocyst
6 differentiation or adaptation to N-stress (**Table S6**). Additionally, using primer
7 extension and northern blot analysis, we have confirmed several TSS in complex
8 promoter regions (**Fig. 2**) or corresponding to asRNAs or ncRNAs (**Fig. 5**). Although
9 a certain percentage of false positives cannot be excluded, when considering
10 potential -10 elements, 9,885 TSS remain in the dataset at a threshold of +3.0 (**Table**
11 **S1**). Approximately one third of all TSS were located upstream of an annotated gene,
12 another third were found within annotated genes, while the remaining TSS were on
13 the reverse complementary strand of 2,412 genes, suggesting antisense transcription
14 to 39% of all genes. This number seems high but is consistent with observations for
15 several other bacteria (31). A total of 1,414 TSS located in the intergenic regions
16 more than 200 nt away from any annotated gene indicated a high number of
17 ncRNAs, although some of these nTSS drive the transcription of mRNAs with very
18 long leaders and therefore are functional gTSS (see comments in **Table S1**).

19 Our data provide unprecedented insight into the complexity of the primary
20 transcriptome of *Anabaena* 7120 under standard growth conditions and at an early
21 step of the heterocyst differentiation process. The use of the *hetR* strain, unable to
22 start the transcriptional program leading to heterocyst differentiation, allowed us to
23 separate transcriptional changes specifically related to this developmental process
24 (DIF) from other N-stress responses that are still observed in the *hetR* mutant (DEF),
25 thus likely unrelated to heterocyst differentiation but rather involved in other aspects
26 of the adaptation to N stress. We thus defined sets of specifically regulated

1 promoters belonging to the DEF and DIF categories. The use of TSS included in
2 these two categories defined both the NtcA and the HetR regulons. As exemplified by
3 the cases of *hetR* (**Fig. 1A**), *nblA* or *alr3808* (**Fig. 2**), the use of TEX-treated samples
4 allowed the identification of multiple TSS in a given promoter region. Complex
5 promoter regions with several TSS are commonly found in genes involved in
6 heterocyst differentiation and patterning, probably due to differential TSS use in the
7 two cell types of the filament (e.g., TSS for *ntcA*, *hetR*, *devB*, *hetC* in **Table S2**).

8 Over the last decades, genes involved in heterocyst differentiation and N₂
9 fixation were primarily identified by mutagenesis and screening of strains unable to
10 grow in the absence of combined nitrogen (2, 30). Here, comparison of the wild-type
11 to the *hetR* transcriptome upon N step-down yielded the DIF category, i.e. the HetR
12 regulon, which now can be analyzed in search of HetR-dependent TSS
13 corresponding to new and previously unknown genes potentially involved in the
14 differentiation process. This regulon includes, for instance, genes related to cell wall
15 synthesis and/or remodeling (**Fig. 3**), a key aspect of heterocyst differentiation. The
16 DIF category also provides a valuable dataset to identify sequence motifs potentially
17 involved in heterocyst-specific expression. Indeed, we identified the DIF+ motif
18 common to many heterocyst-specifically expressed promoters. It consists of a short
19 palindrome 5' TCCGGA, centered at or close to position -35, suggesting it might be
20 recognized by a specific sigma factor rather than serve as a HetR binding site.

21 NtcA-mediated regulation is operated by binding to a consensus sequence,
22 GTAN₈TAC, first described for strongly regulated promoters in *Synechococcus* (5).
23 Promoters that are directly activated by NtcA contain an NtcA binding site that, in
24 most cases, is centered close to position -41.5 with respect to the TSS (although
25 some NtcA binding sites further upstream are also described). In such promoters,
26 NtcA activates transcription in a manner that resembles CAP-mediated regulation at

1 Class II promoters. On the other hand, NtcA-mediated repression is operated by
2 interaction with NtcA binding sites at a position that makes binding incompatible with
3 the normal operation of the promoter. In the case of *gifA* from *Anabaena*, one NtcA
4 binding site is centered at position -28.5 (21). The DEF category defined here
5 identifies transcriptional responses that in many cases were directly regulated by
6 NtcA as deduced from the identification of NtcA binding sites located in positions
7 compatible with transcriptional activation (DEF+) or repression (DEF-) with respect to
8 the corresponding TSS (**Tables S8, S9** and **Fig. 4**). Comparison to a computational
9 prediction (32) of the NtcA regulon (**Table S10**) revealed that many of the TSS in the
10 DEF category correspond to previously unknown NtcA-regulated promoters, thereby
11 expanding the known NtcA regulon. The absence of NtcA binding sites in the
12 promoters for several other TSS in the DEF categories (such as TSS1 for *nblA*; **Fig.**
13 **2B and 4A**) suggests their expression is either not directly regulated by NtcA or
14 operated by binding to different positions.

15 Finally, as observed in other cyanobacteria (15, 33, 34), our dataset indicates
16 the abundant transcription of antisense and ncRNAs (e.g., ncRNA T1 in **Fig. 1B**,
17 associated with >600,000 reads). As previously described for NsiR1, a short ncRNA,
18 whose expression is induced specifically in proheterocysts upon N step-down and
19 belongs into the NtcA/HetR regulon (25), expression of some of these transcripts
20 (**Fig. 5**) is regulated by N availability, suggesting that antisense and non-coding
21 transcripts might be involved in the regulation of nitrogen assimilation and heterocyst
22 differentiation. The annotated primary transcriptome of *Anabaena* 7120 during the
23 transition from ammonium utilization to N₂ fixation will greatly facilitate the use of this
24 organism as a model for prokaryotic cell differentiation and N₂ fixation in an oxygenic
25 phototroph.

26

1 **METHODS**

2 Full protocols are available in SI Materials and Methods.

3

4 **Growth conditions.** Cultures of *Anabaena* 7120 WT, *hetR* mutant 216 (8) and *ntcA*
5 mutant CSE2 (6) were bubbled with an air/CO₂ mixture (1% v/v) and grown
6 photoautotrophically at 30°C in BG110C medium lacking NaNO₃ but containing 6 mM
7 NH₄Cl, 10 mM NaHCO₃, and 12 mM N-tris (hydroxymethyl) methyl-2-
8 aminoethanesulfonic acid-NaOH buffer (pH 7.5). Four RNA samples were isolated for
9 dRNA-seq analysis from cells taken at T = 0 h (WT-0 and *hetR*-0) and T = 8 h (WT-8
10 and *hetR*-8) after removing all combined nitrogen from the media.

11

12 **Preparation and analysis of RNA.** Total RNA was isolated using hot phenol (35)
13 with modifications. Northern blot hybridization and primer extension analysis of 5'
14 ends was performed as described (13, 34, 36).

15 **Deep transcriptome sequencing.** The cDNA libraries were prepared by vertis
16 Biotechnologie, Germany (<http://www.vertisbiotech.com/>) after enrichment for primary
17 transcripts by treatment with TEX (Epicentre) and analyzed on an Illumina sequencer
18 by Beckman Coulter Genomics, Danvers, MA as previously described (16). Based on
19 tetranucleotide tags (**Table S11**), 5.153.094, 4.690.212, 6.398.708 and 5.497.219
20 from these sequence reads were assigned to the WT-0, the WT-8, the *hetR*-0 and
21 *hetR*-8 populations, respectively, and matched against the sequences of the
22 chromosome or plasmids of *Anabaena* 7120.

23 **Computational methods.** Reads <18 nt and those with blast hits to the ribosomal
24 clusters were filtered out. Remaining reads were mapped to the genome using the
25 *segemehl* algorithm (37), with default parameters. Reads were pooled from the four
26 samples and their 5' ends were binned within a 5-nt section. The position within the

1 window where the most reads began was considered to be the initial TSS. Because
2 we noticed a few cases of initiation of transcription from a broader window, this
3 dataset was clustered to allow the combination of initial TSS, which were not further
4 than 5 nt apart. The position within this window where the greatest number of reads
5 began was considered a TSS when a minimum of 50 reads was associated with it.

6 For ratio calculation, the number of reads for the four samples were
7 normalized (19), and single pseudocounts were added to make the calculation of
8 ratios possible for all TSS. The resulting ratios were classified and filtered into DEF
9 and DIF categories of regulated promoters. Possible -10 elements were searched 6-8
10 nt upstream of all putative TSS and scored according to a PSSM derived from this
11 dataset (**Table S1, Fig. S6**). To construct a PSSM for the NtcA binding site, all
12 promoter regions for the 129 TSS in the DEF+ category with at least 8-fold change
13 were aligned. From these, 81 possessed an element matching at least 4 nt of the
14 GTAN₈TAC motif centered at position 22/23 upstream of the -10 element. These
15 elements were used together with six additional experimentally defined sites (**Table**
16 **S2**) in the construction of the matrix. For the DIF+ motif, the regions -44 to -25 of all
17 mapped TSS were searched in two distinct datasets: One consisted of all 209 TSS in
18 the DIF+ class and the second of the remaining 13,496 putative TSS. Statistical
19 significance of an enrichment for the DIF+ motif in the DIF+ class was tested in a
20 Pearson's chi-squared test.

21

22 **ACKNOWLEDGMENTS**

23 This work was supported by the DFG program “Sensory and regulatory RNAs in
24 Prokaryotes” SPP1258 (grant HE 2544 4-2), by the BMBF grant 0313921 (WRH) and
25 by the Ministerio de Ciencia e Innovación grants BFU2007-60651 (AV) and
26 BFU2010-14821 (AMP) co-financed by FEDER.

1 **References**

2

- 3 1. Flores E, Herrero A (2010) Compartmentalized function through cell
4 differentiation in filamentous cyanobacteria. *Nat Rev Microbiol* 8:39-50.
- 5 2. Kumar K, Mella-Herrera RA, Golden JW (2010) Cyanobacterial heterocysts.
6 *Cold Spring Harbor Persp Biol* 2:1-19.
- 7 3. Herrero A, Muro-Pastor AM, Flores E (2001) Nitrogen control in cyanobacteria.
8 *J Bacteriol* 183:411-425.
- 9 4. Herrero A, Muro-Pastor AM, Valladares A, Flores E (2004) Cellular
10 differentiation and the NtcA transcription factor in filamentous cyanobacteria.
11 *FEMS Microbiol Rev* 28:469-487.
- 12 5. Luque I, Forchhammer K (2008) Nitrogen assimilation and C/N balance
13 sensing. *The Cyanobacteria: Molecular Biology, Genomics and Evolution*, eds
14 Herrero A & Flores E (Caister Academic Press, Hethersett), pp 335-382.
- 15 6. Frías JE, Flores E, Herrero A (1994) Requirement of the regulatory protein NtcA
16 for the expression of nitrogen assimilation and heterocyst development genes in
17 the cyanobacterium *Anabaena* sp. PCC7120. *Mol Microbiol* 14:823-832.
- 18 7. Wei TF, Ramasubramanian TS, Golden JW (1994) *Anabaena* sp. strain PCC
19 7120 *ntcA* gene required for growth on nitrate and heterocyst development. *J*
20 *Bacteriol* 176:4473-4482.
- 21 8. Buikema WJ, Haselkorn R (1991) Characterization of a gene controlling
22 heterocyst differentiation in the cyanobacterium *Anabaena* 7120. *Genes &*
23 *Development* 5:321-330.
- 24 9. Black TA, Cai Y, Wolk CP (1993) Spatial expression and autoregulation of *hetR*,
25 a gene involved in the control of heterocyst development in *Anabaena*. *Mol*
26 *Microbiol* 9:77-84.

- 1 10. Huang X, Dong Y, Zhao J (2004) HetR homodimer is a DNA-binding protein
2 required for heterocyst differentiation, and the DNA-binding activity is inhibited
3 by PatS. *Proc Natl Acad Sci USA* 101:4848-4853.
- 4 11. Kim YC, *et al.* (2011) Structure of HetR - a novel transcription factor required for
5 heterocyst differentiation in cyanobacteria. *Proc Natl Acad Sci USA* 108:10109-
6 10114.
- 7 12. Higa KC, Callahan SM (2010) Ectopic expression of *hetP* can partially bypass
8 the need for *hetR* in heterocyst differentiation by *Anabaena* sp. strain PCC
9 7120. *Mol Microbiol* 77:562-574.
- 10 13. Muro-Pastor AM, Valladares A, Flores E, Herrero A (2002) Mutual dependence
11 of the expression of the cell differentiation regulatory protein HetR and the
12 global nitrogen regulator NtcA during heterocyst development. *Mol Microbiol*
13 44:1377-1385.
- 14 14. Flaherty BL, Van Nieuwerburgh F, Head SR, Golden JW (2011) Directional
15 RNA deep sequencing sheds new light on the transcriptional response of
16 *Anabaena* sp. strain PCC 7120 to combined-nitrogen deprivation. *BMC*
17 *Genomics* 12:332.
- 18 15. Mitschke J, *et al.* (2011) An experimentally anchored map of transcriptional start
19 sites in the model cyanobacterium *Synechocystis* sp. PCC6803. *Proc Natl Acad*
20 *Sci USA* 108:2124-2129.
- 21 16. Sharma CM, *et al.* (2010) The primary transcriptome of the major human
22 pathogen, *Helicobacter pylori*. *Nature* 464:250-255.
- 23 17. Buikema WJ, Haselkorn R (2001) Expression of the *Anabaena hetR* gene from
24 a copper-regulated promoter leads to heterocyst differentiation under repressing
25 conditions. *Proc Natl Acad Sci USA* 98:2729-2734.

- 1 18. Rajagopalan R, Callahan SM (2010) Temporal and spatial regulation of the four
2 transcription start sites of *hetR* from *Anabaena* sp. strain PCC 7120. *J Bacteriol*
3 192:1088-1096.
- 4 19. Robinson MD, Oshlack A (2010) A scaling normalization method for differential
5 expression analysis of RNA-seq data. *Genome Biol* 11:R25.
- 6 20. Frías JE, Flores E, Herrero A (1997) Nitrate assimilation gene cluster from the
7 heterocyst-forming cyanobacterium *Anabaena* sp. strain PCC 7120. *J Bacteriol*
8 179:477-486.
- 9 21. Galmozzi CV, Saelices L, Florencio FJ, Muro-Pastor MI (2010)
10 Posttranscriptional regulation of glutamine synthetase in the filamentous
11 cyanobacterium *Anabaena* sp. PCC 7120: differential expression between
12 vegetative cells and heterocysts. *J Bacteriol* 192:4701-4711.
- 13 22. Baier K, Lehmann H, Stephan DP, Lockau W (2004) NblA is essential for
14 phycobilisome degradation in *Anabaena* sp. strain PCC 7120 but not for
15 development of functional heterocysts. *Microbiology* 150:2739-2749.
- 16 23. Ehira S, Ohmori M (2006) NrrA, a nitrogen-responsive response regulator
17 facilitates heterocyst development in the cyanobacterium *Anabaena* sp. strain
18 PCC 7120. *Mol Microbiol* 59:1692-1703.
- 19 24. Ow SY, *et al.* (2008) Quantitative shotgun proteomics of enriched heterocysts
20 from *Nostoc* sp. PCC 7120 using 8-plex isobaric peptide tags. *J Proteome Res*
21 7:1615-1628.
- 22 25. Ionescu D, Voss B, Oren A, Hess WR, Muro-Pastor AM (2010) Heterocyst-
23 specific transcription of NsiR1, a non-coding RNA encoded in a tandem array of
24 direct repeats in cyanobacteria. *J Mol Biol* 398:177-188.
- 25 26. Valladares A, *et al.* (2011) Specific role of the cyanobacterial PipX factor in the
26 heterocysts of *Anabaena* sp. strain PCC 7120. *J Bacteriol* 193:1172-1182.

- 1 27. Muro-Pastor AM, Olmedo-Verd E, Flores E (2006) All4312, an NtcA-regulated
2 two-component response regulator in *Anabaena* sp. strain PCC 7120. *FEMS*
3 *Microbiol Lett* 256:171-177.
- 4 28. Valladares A, Muro-Pastor AM, Herrero A, Flores E (2004) The NtcA-dependent
5 P1 promoter is utilized for *glnA* expression in N₂-fixing heterocysts of *Anabaena*
6 sp. strain PCC 7120. *J Bacteriol* 186:7337-7343.
- 7 29. Ramasubramanian TS, Wei TF, Golden JW (1994) Two *Anabaena* sp. strain
8 PCC 7120 DNA-binding factors interact with vegetative cell- and heterocyst-
9 specific genes. *J Bacteriol* 176:1214-1223.
- 10 30. Lechno-Yossef S, Fan Q, Wojciuch E, Wolk CP (2011) Identification of ten
11 *Anabaena* sp. genes that, under aerobic conditions, are required for growth on
12 dinitrogen but not for growth on fixed nitrogen. *J Bacteriol* 193:3482-3489.
- 13 31. Georg J, Hess WR (2011) cis-antisense RNA, another level of gene regulation
14 in bacteria. *Microbiol Mol Biol Rev* 75:286-300.
- 15 32. Novichkov PS, *et al.* (2010) RegPrecise: a database of curated genomic
16 inferences of transcriptional regulatory interactions in prokaryotes. *Nucleic*
17 *Acids Res* 38:D111-118.
- 18 33. Georg J, *et al.* (2009) Evidence for a major role of antisense RNAs in
19 cyanobacterial gene regulation. *Mol Syst Biol* 5:305.
- 20 34. Steglich C, *et al.* (2008) The challenge of regulation in a minimal
21 photoautotroph: non-coding RNAs in *Prochlorococcus*. *PLoS Genet*
22 4:e1000173.
- 23 35. Mohamed A, Jansson C (1989) Influence of light on accumulation of
24 photosynthesis-specific transcripts in the cyanobacterium *Synechocystis* 6803.
25 *Plant Mol Biol* 13:693-700.

- 1 36. Muro-Pastor AM, Valladares A, Flores E, Herrero A (1999) The *hetC* gene is a
2 direct target of the NtcA transcriptional regulator in cyanobacterial heterocyst
3 development. *J Bacteriol* 181:6664-6669.
- 4 37. Hoffmann S, *et al.* (2009) Fast mapping of short sequences with mismatches,
5 insertions and deletions using index structures. *PLoS Comput Biol* 5:e1000502.
- 6 38. Vioque A (1992) Analysis of the gene encoding the RNA subunit of
7 ribonuclease P from cyanobacteria. *Nucleic Acids Res* 20:6331-6337.
- 8 39. Aldea MR, Mella-Herrera RA, Golden JW (2007) Sigma factor genes *sigC*, *sigE*,
9 and *sigG* are upregulated in heterocysts of the cyanobacterium *Anabaena* sp.
10 strain PCC 7120. *J Bacteriol* 189:8392-8396.
- 11 40. Muro-Pastor AM, Flores E, Herrero A (2009) NtcA-regulated heterocyst
12 differentiation genes *hetC* and *devB* from *Anabaena* sp. strain PCC 7120 exhibit
13 a similar tandem promoter arrangement. *J Bacteriol* 191:5765-5774.
- 14

1 **FIGURE LEGENDS**

2

3 **Figure 1. Genome-wide identification of transcriptional start sites (TSS) in**
4 ***Anabaena* 7120.** (A) Differential RNA-seq of TEX-treated samples identifies single
5 TSS in complex promoter regions as exemplified by the *hetR* gene. The total number
6 of reads mapped to each 5' end is indicated for each of the four previously described
7 TSS. (B) Distribution of 3401 TSS with ≥ 300 reads each along a linear plot of the
8 *Anabaena* 7120 chromosome. TSS mapped for the forward strand are plotted above
9 the x-axis, and for the reverse strand below. The number of sequence reads is given
10 on the y-axis (logarithmic scaling). The location of each TSS according to **Fig. S1A**
11 served for classification as gTSS (blue), nTSS (green), aTSS (red) or iTSS (grey).
12 Selected TSS for each of the four classes are annotated.

13

14 **Figure 2. Analysis of genes with multiple TSS identified by dRNA-seq.** RNA was
15 isolated from ammonium-grown cells (lanes labeled 0) or from ammonium-grown
16 cells incubated in the absence of combined nitrogen for the number of hours
17 indicated. (A) Graphical representation of reads mapped to the promoter of *nbIA*. The
18 histograms correspond to the WT-0 (red), WT-8 (green), *hetR*-0 (black) and *hetR*-8
19 (blue) samples. (B) Primer extension analysis of the *nbIA* mRNA in WT and mutant
20 strains CSE2 (*ntcA*) and 216 (*hetR*). Samples contained 20 μg of RNA. The
21 oligonucleotide used was complementary to positions +5 to -18 with respect to the
22 translational start of *nbIA*. The 5' ends identified by dRNA-seq are numbered 1-5 (for
23 positions, see **Table S1**). (C) Graphical representation of reads mapped to the
24 promoter of *alr3808*. (D) Northern blot analysis of the *alr3808* mRNA in *Anabaena*
25 7120 and mutant strain 216 (*hetR*). Samples contained 10 μg of RNA. The probe
26 used was an internal fragment of *alr3808*. *mnpB* (38) was used as a loading control.

1

2 **Figure 3. Occurrence of a palindrome, 5' TCCGGA, in promoters of the DIF+**
 3 **category.** (A) Alignment of the heterocyst-specific promoters for NsiR1 (25), the *hetR*
 4 TSS3 (18), *sigC* (39) and *hetC* (36, 40) with selected promoters in the DIF+ category
 5 (fold change ≥ 8) containing TCCGGA around position -35 (one mismatch allowed).
 6 (B) Cell-specific transcription from the wild-type promoter of NsiR1 (P6; upper
 7 panels) or a mutated version of P6 carrying GAATTC instead of TCCGGA (lower
 8 panels). Images corresponding to red autofluorescence (left panel) and GFP
 9 fluorescence (right panel) are shown. White triangles point to proheterocysts. Scale
 10 bars, 10 μm .

11

12 **Figure 4. NtcA-activated and repressed promoters.** (A) Promoter regions of 20
 13 TSS in the DEF+ category with the highest fold change. Possible -10 elements are
 14 highlighted in blue, nucleotides matching the consensus for NtcA binding sites in red,
 15 and the underlined nucleotide in each sequence is the TSS. (B) Nucleotide
 16 frequencies derived for positions 1 to 14 of 87 putative NtcA binding sites. (C)
 17 Corresponding Weblogo. (D) Position of NtcA binding sites identified in a sliding
 18 window approach using the PSSM along the promoters in the DEF+ category (fold
 19 change ≥ 8). The bars indicate the first nucleotide of a putative NtcA binding site and
 20 its position with regard to the TSS (E) Putative NtcA binding sites identified at
 21 repression-compatible positions around TSS in the DEF- category.

22

23 **Figure 5. Experimental verification of newly identified transcripts classified as**
 24 **antisense or non-coding.** RNA was isolated from ammonium-grown cells (lanes
 25 labeled 0) or from cells incubated in the absence of combined nitrogen for the
 26 number of hours indicated. (A) Primer extension analysis of the *all3278* asRNA in

1 *Anabaena* 7120 WT and mutant strains CSE2 (*ntcA*) and 216 (*hetR*). Samples
2 contained 20 µg of RNA. (B) Northern blot analysis of ncRNAs NsiR2 (upper panel)
3 and NsiR3 (middle panel) in *Anabaena* 7120 WT and mutant strains CSE2 (*ntcA*)
4 and 216 (*hetR*). The samples contained 10 µg of RNA. The *trnL-UAA* transcript
5 (lower panel) was used as a loading control. (C) Predicted secondary structure of
6 NsiR3.

Table 1. Overview on the number and types of putative TSS mapped for the chromosome (chr) and plasmids alpha, beta, gamma, delta, epsilon and zeta.

| | chr | alpha | beta | gamma | delta | epsilon | zeta | Total |
|--------------------|------------|--------------|-------------|--------------|--------------|----------------|-------------|--------------|
| length (nt) | 6,413,771 | 408,101 | 186,614 | 101,965 | 55,414 | 40,340 | 5,584 | - |
| # of genes | 5430 | 386 | 186 | 90 | 85 | 31 | 5 | |
| gTSS | 3955 | 145 | 41 | 24 | 14 | 6 | 1 | 4186 |
| aTSS | 3854 | 188 | 73 | 34 | 13 | 8 | 2 | 4172 |
| iTSS | 3722 | 113 | 50 | 18 | 21 | 6 | 3 | 3933 |
| nTSS | 1266 | 88 | 15 | 22 | 8 | 7 | 8 | 1414 |
| Total | 12797 | 534 | 179 | 98 | 56 | 27 | 14 | 13705 |

FIGURE 1

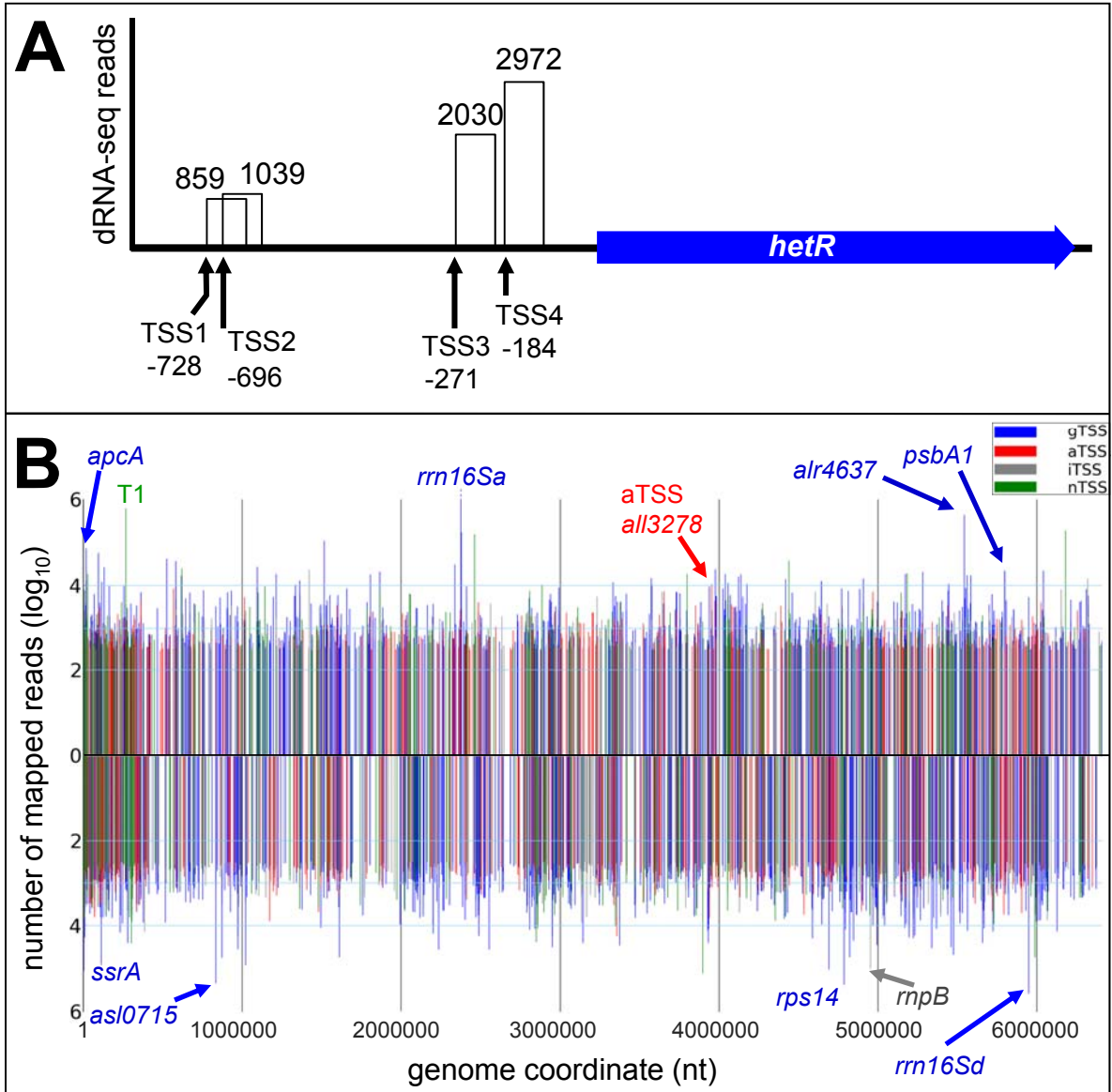


FIGURE 2

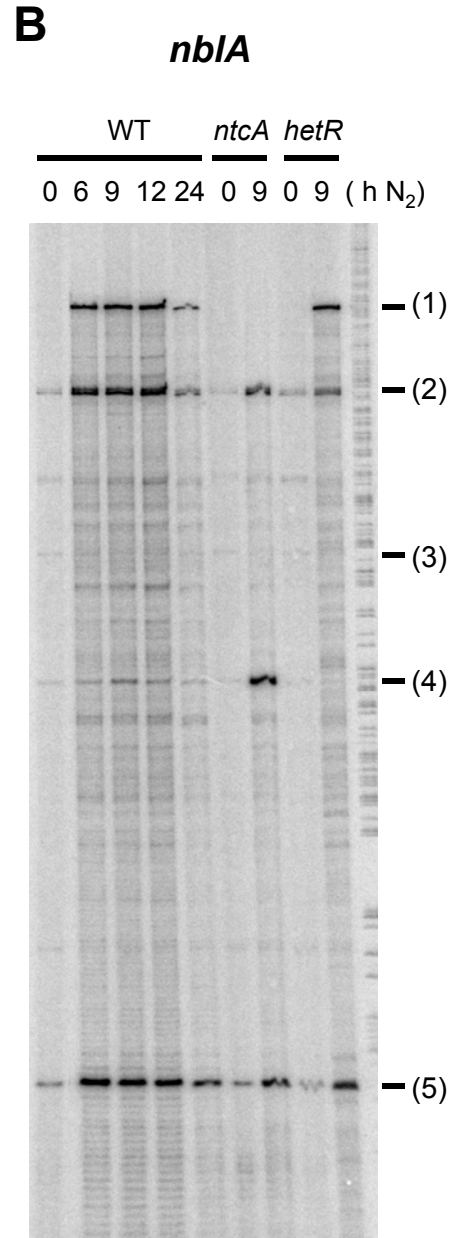
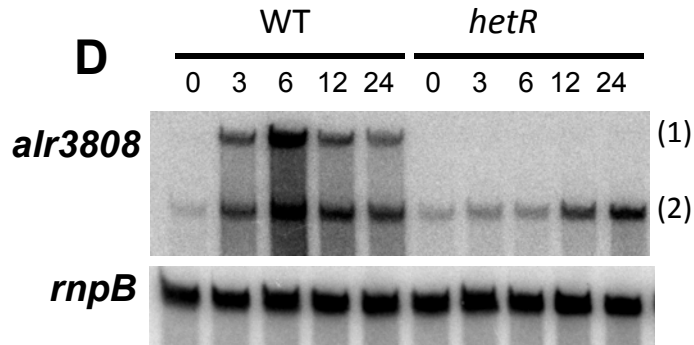
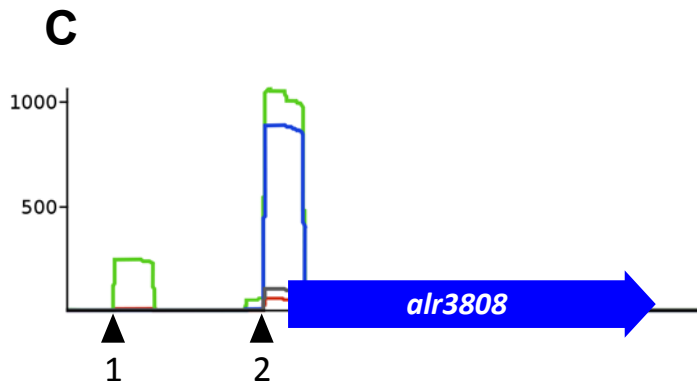
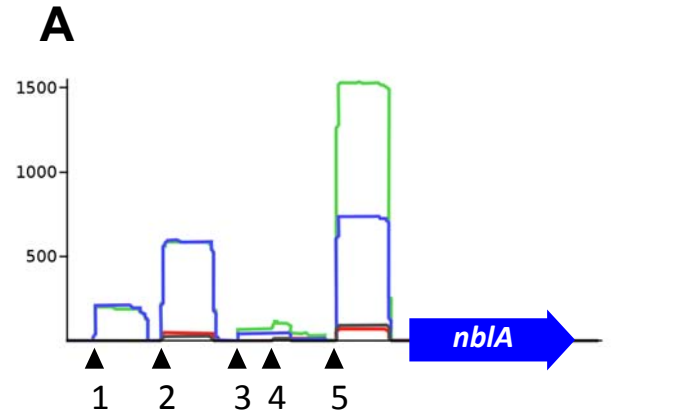


FIGURE 3

| | | | | | |
|---|----------|---------------|--------|----------------------------------|----------------------------------|
| A | 4271758r | TTGATGCAATTTA | TCCGGA | AGACTGTAATTCAAAATAGAACAAATTAATTG | NsiR1 (TSS12) Table S2 |
| | 2821366f | AGAGCAGATAAGT | TCCGGA | TAATAGGGAAAGTCCTTGTAGGTTACTTATTA | hetR (TSS3) Table S2 |
| | 2022661r | AATTAGAAAATCG | TCAGGA | AATTACTTATATACCCATGTAGATGTGACTA | sigC Table S2 |
| | 3427771f | AAAAAATAATTT | TCCTGA | TGTTTTAAGAAAATCTGTTGTTATAAATT | hetC (TSS2) Table S2 |
| | 3453831f | AGAATTAGGTTTA | TCCTGA | AAGAGTAAAAAAAATCCGAATATCCTAAATTA | hepA HEP island Table S6 |
| | 3449710f | AGCATTATTAGAG | TCCGGA | GAAATCACTTGGGATTATGAAAATATTCGTA | alr2833 HEP island Table S6 |
| | 3457231f | GTTCCGATGTTCA | TACGGA | AAGCACCACAATTTAGCCGTAAGTATGTTTA | alr2837 HEP island Table S6 |
| | 580704f | AGAGGTATTATTG | TCCGAA | TATTTGTCTTTCACCTGCGAAAAAATTATA | asr0485 (<i>pipX</i>) Table S2 |
| | 4601709f | TTCTGATAATTTT | TCCTGA | GAACACCATTATTTACAAGTAGAGTGTGATG | alr3808 Figure 2 |
| | 3569154 | GTTGTGCGCCCTT | CCCGGA | TTGTAGCGATCCGAGTAGAACCTGTTTTTA | alr2933 peptidoglycan turnover |
| | 5742628r | TCCCACAGCACCT | TCCGGA | AAATTAAAAAACGCCGAAAATATTGCA | all4822 cell envelope biogenesis |
| | 5953754f | TTACACAATTTTA | TCGGA | TATATACCGCTTCAGGTGGAATAAATTTCTTA | alr4984 peptidoglycan-binding |
| | 5954131f | AAGTAGACGTTAT | TCCGGA | ATAATTGATTTTTTTCTGTCAAAGGAATTA | alr4984 peptidoglycan-binding |
| | 518833r | ATCCCAATTCAA | TCCGGA | ATATTCTATAATCTGGAGAATAATAAATTATC | all0438 ser/thr kinase |
| | 3072950r | TTAAGAAAAGTT | TCCGGA | TTTGAGTCCCAATGAGTGTCTAAGTTGTTA | all2571 unknown protein |
| | 2881051r | AGGTATTTCTGTA | TCCGGA | TAATTAACCTGCCAAGTGTAAATCTTAACA | as12397 unknown protein |

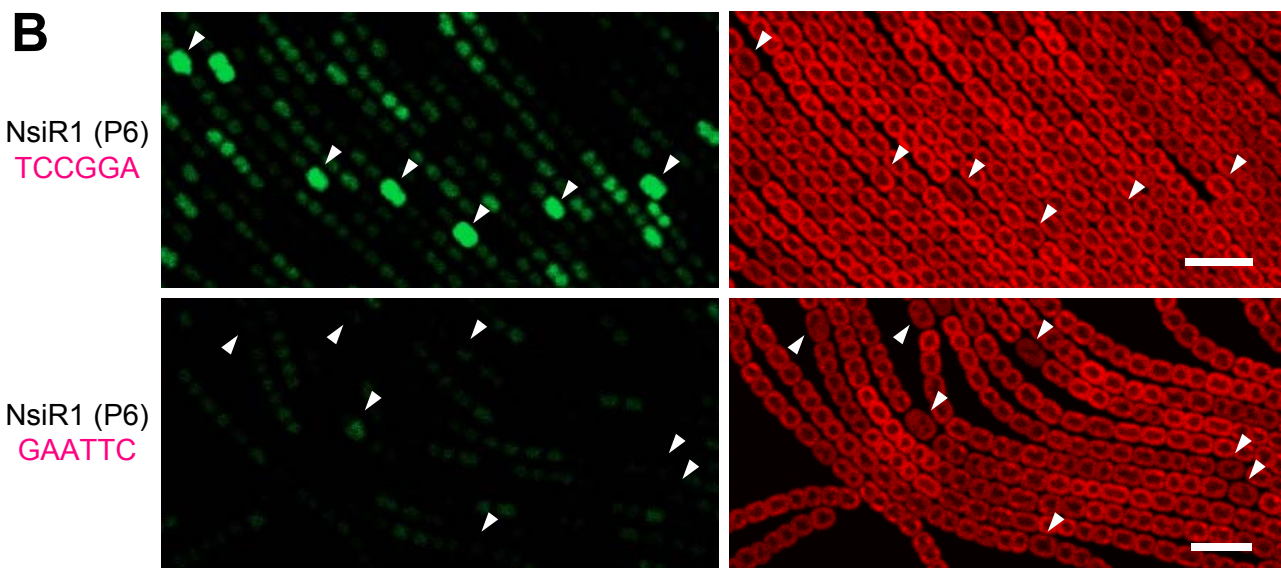


FIGURE 4

A

```

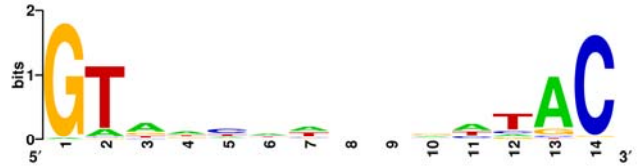
2058877f AGTTTAT GTAACCTATAAGAC ATTTTATTTGATACCTCATACTC TAAAAT CAAGTA nTSS
2460393r TAAAAA GTAACTCTTGATAC ATACGCTTATGAAAAACGCATA TACCAT TGAAAAA gtSS asl2052
580293f TAAAAA GTAGCAATGCAGAC TGTTGTTAGGAACAGTTATTAG GAGAAT GCGCCTG gtSS asr0485 (pipX) Table S2
1273249r TCTTTTG GGTACAAGATATAC AAAATAATATTGAGGAATTAGGC TATCTT CATATCT gtSS all1087
5547631f GTTTTT GTTGCGTGCTAGAC ATAACCAGACGGGTGTTTTGATC CAAACT CCTGTA atSS all14644
2059119f TTATTTT GTATTTAAACGGGAC AGTTCCTACTTATCTAGTTAAGT TTAAAT AACAACTA gtSS alr1713
2837125f GTAGATA GATATCCACATAC GGAAGTGTCACTCTGATACTGG CAGGCT AAATTA gtSS alr2355
5731963f GTTTGTT GGCGCAACGGCTAC AGTTTGTGGCGAGACAGGG GATGAT GGATTAG atSS all14813
4907756f GCAAAC GAATGTTTGATAC GGCAGGATGTGCAGTTTCTCT TACCTT GAGCAAG gtSS alr4077
1693413r AAAAAA GTAATCAGCCTGAC AGAACTATCGTCTGATTAGGAG TATAAA GTGATCA gtSS all1432 (hesA) Table S2
3953418f TGAGTTA GTCGCTAAAGCTAC ATTTTGCTAACAGTATCCGACT TATTAT GAGATTTA atSS all13278 Figure 5I
2400767r GTTGCTC GTATATTTCAACAC GAATTTGATCATTAGATGGTG TACTGT TTATAGA gtSS all2006 Table S6
519953f ACATAAC GTGTTTTCAGTTAC AGTTATGCCAGATGCAATTAAGC CACAAT GTTGATTA gtSS alr0440
105428r CATTATG GTATGAAATAGTAC AGTTTAAAAATTAGTGTTCGCGT CATCAT TACGAG atSS all17614
1657401r GAGAGTC GTAGCAATAACACAC TAAAACCTCTGGAACAGTAGGT TAGGCT TGCCTTA gtSS all1395
3346518f ATAAACT GATAGTTATATAC TGTTCTCAGAAACGAAAACTA TATATT GAGCATA nTSS
5248514r TGTTTT GCGATCGGCGATAC AATTTACACGGGGCAAAAGCTG GAATAT GAAGGA itSS all14379
5167792r GGCTAGA GTACAAAAGCTAC AAAACCTTGGGCATGGGCTTGT TACTTT GAAATTC gtSS all14312 (nrrA) Table S2
5407066f CTCAGCAATTTGTTCAACCTGAGCATTTCACCATTTGCAACTTGA TACAAA TATTTTTA gtSS asr4517 (nblA) Table S6
2793917r CTTCTCAACTGCTCATAAGACAGATACGGTTAAAAAAGTTGC AATTCT CATAAGT gtSS all12319 (glnB) Table S6

```

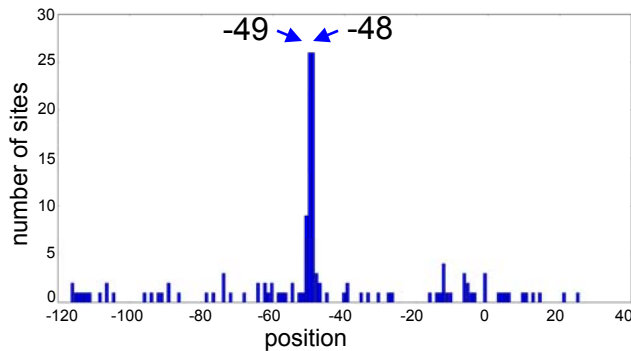
B

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
|---|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| A | 2 | 10 | 47 | 35 | 16 | 36 | 34 | 28 | 24 | 29 | 39 | 11 | 70 | 1 |
| C | 0 | 2 | 8 | 7 | 40 | 14 | 11 | 17 | 13 | 10 | 17 | 12 | 4 | 83 |
| G | 85 | 2 | 17 | 21 | 8 | 12 | 8 | 16 | 21 | 32 | 4 | 10 | 10 | 2 |
| T | 0 | 73 | 15 | 24 | 23 | 25 | 34 | 26 | 29 | 16 | 27 | 54 | 3 | 1 |

C



D



E

```

2828946f GGATCAATACGATTTTCTATCTAGAAG GAAAAT ATAGAGAC GAAATCAGCCCTAAACGGCACTAAAGATGAA gtSS alr2346
3359637f TTTTGATTTTTCATATACAATTTTGT TATTTT TAGATTAAGATGTTATTTACCTAAATTTAATCAACGTGT gtSS asr2763
5441031f TTACTGGAAAAATTAGCGGTTTGACACA TAGAAT TATTAGAGAC GGTAATAATTTGTAAACTAGATTGACAAT gtSS alr4548 (psbD)
5046301r TTGAGATTTTAGATATGGAAATTAATTACAAT CTAATAATCGAAAAATAAAGTTTCCCTGACAAAGCCCAGG gtSS all14203 (rpl5)
2809313r CCGTAGCATAAGATACAGAATCTTGC TATATT AAATGTGTGAAGGTCAAATCCAATTAATCACTAGGA gtSS asl2329 (gifA)
3433648r GTCTTTTAAACACGCCCAATCACTGC CATGAT GGATTAGTTCATAGTGTGTGTTGTTGTTGTGATAGGCGG nTSS
4785763f CCAGGCAATCATTGTCATCAGCCATGA GAAAAT AGCCTTAGCGTAGCCTTGAATAA TAGAGCGATATCAT gtSS alr3968
1577153f CTSTAACATACACTACGAAACTTATGC TATGTT AGGAAGAACCAGACATAAAGCAGAAAAATTAAGAGGTTAA gtSS asr1328
2807328f CTTTTGTGCAGATGTGAAAGAAAGGT TAATAT TACCTGTAATCCAGACGTTCTGTACAAAAGCTACAAA gtSS alr2328 (glnA)
1785466f CAAAGAATAACTTATGCCATTTCTTGA TATATT GTGAGACAAGTTACAAATTACGTGGTGTGCAATTTTTTC gtSS rbcL

```

FIGURE 5

