# Widespread divergent transcription from bacterial and archaeal promoters is a consequence of DNA sequence symmetry

Warman, Emily; Forrest, David; Guest, Thomas; Haycocks, James; Wade, Joseph T.; Grainger, David

*Document Version*
Peer reviewed version

# Widespread divergent transcription from bacterial and archaeal promoters is a consequence of DNA sequence symmetry

**Emily A. Warman[1], David Forrest[1], Thomas Guest[1], James J.R.J. Haycocks[1],**

**Joseph T. Wade[2,3], David C. Grainger[1]***

[1]Institute for Microbiology and Infection, School of Biosciences, University of Birmingham,

Edgbaston, Birmingham, B15 2TT, UK

[2]Wadsworth Centre, New York State Department of Health, Albany, NY, 12208, USA

[3]Department of Biomedical Sciences, University at Albany, Albany, NY, 12201, USA

*for correspondence, d.grainger@bham.ac.uk Tel: +44 (0)121 4145437

1 **ABSTRACT**

2 **Transcription initiates at promoters, DNA regions recognised by a DNA-dependent RNA**

3 **polymerase. We previously identified horizontally acquired *Escherichia coli* promoters where the**

4 **direction of transcription was unclear. Here, we show that more than half of these promoters are**

5 **bidirectional. Using genome-scale approaches, we demonstrate that 19% of all transcription start**

6 **sites detected in *E. coli* are associated with a bidirectional promoter. Bidirectional promoters are**

7 **similarly common in diverse bacteria and archaea and have inherent symmetry: specific bases**

8 **required for transcription initiation are reciprocally co-located on opposite DNA strands.**

9 **Bidirectional promoters enable co-regulation of divergent genes and are enriched in both**

10 **intergenic and horizontally acquired regions. Divergent transcription is conserved among**

11 **bacteria, archaea and eukaryotes, but the underlying mechanisms for bidirectionality are**

12 **different.**

13

14 **INTRODUCTION**

15 Promoters are sections of duplex DNA that interact with RNA polymerase (RNAP) to stimulate

16 transcription initiation[1]. In most organisms, promoters consist of ordered core elements with distinct

17 roles[2,3]. For instance, the bacterial -10 element (consensus 5'-TATAAT-3') is usually indispensable and

18 interacts with the housekeeping RNAP $\sigma^{70}$ subunit ($\sigma^{A}$ in some bacteria). The second and sixth positions

19 of -10 elements are most critical; non-template strand bases interact with $\sigma^{70}$ to stabilise DNA

20 unwinding[4,5]. Position one is also important, and defines the upstream boundary of DNA melting[5]. Less

21 conserved ancillary sequences can aid RNAP recruitment. For instance, the -35 element (consensus

22 5'-TTGACA-3') also contacts $\sigma^{70}$. Following initiation, $\sigma^{70}$ is evicted from the elongation complex. In

23 many eukaryotes and archaea, the TATA box functions analogously to the bacterial -10 element; TATA

24 binding protein (TBP) facilitates DNA unwinding and serves as a scaffold for recruiting the

25 transcriptional apparatus[6].

26

27 It has long been assumed that promoter sequences are directional, driving transcription in a single

28 orientation determined by promoter element arrangement[2,7]. This view has been challenged in

29 eukaryotes[8–11]. In addition to driving the production of a canonical sense mRNA, many RNAP II

30 promoters simultaneously stimulate antisense transcription[12]. Permissive chromatin plays a key role;

31 nucleosome-depleted DNA allows fortuitous binding of transcriptional activators that permit divergent

32 transcription[12–15]. Thus, permissive sections of eukaryotic chromatin, not core promoters *per se*, give

33 rise to bidirectionality[9,12,13,16]. The phenomenon is particularly prevalent for recently evolved promoter

34 regions suggesting that, if beneficial, selection can fix mutations that favour unidirectional

35 transcription[16]. In prokaryotic organisms, particularly the bacteria, chromosomes are not folded into

36 structures reminiscent of eukaryotic chromatin[17]. The consensus view remains that transcription from

2

37  bacterial promoter sequences is unidirectional[18]. Here, we show that bidirectional prokaryotic promoter

38  sequences, resulting in divergent transcription, are in fact commonplace. However, the underlying

39  molecular mechanisms are fundamentally different to those in eukaryotes.

40

41  **RESULTS**

42  *Identification of bidirectional promoter sequences in horizontally acquired Escherichia coli genes*

43  Transcription start sites (TSSs) in *Escherichia coli* have been mapped by detecting triphosphorylated

44  RNA 5' ends[19]. These can be assigned to $\sigma^{70}$ binding events identified using ChIP-seq[19]. We noticed

45  that not all $\sigma^{70}$ binding was associated with detectable RNA synthesis. This was particularly evident for

46  horizontally acquired genes silenced by histone-like nucleoid structuring (H-NS) protein (Extended

47  Data Fig. 1). We reasoned that RNAP might initiate transcription, but produce unstable RNAs, at these

48  sites; similar to cryptic unannotated transcripts in eukaryotes[12]. To test this, we transcriptionally fused

49  33 such $\sigma^{70}$ targets, derived from H-NS silenced genes, to *lacZ*. Translation prevents Rho mediated

50  transcription termination. Hence, any RNAs produced should be stabilised and detectable[20]. As

51  transcription orientation cannot be directly inferred from $\sigma^{70}$ ChIP-seq data, DNA sequences were

52  cloned in both directions (Fig. 1a). Over half of the fragments were transcriptionally active ($\geq$ 2-fold

53  above the background control) regardless of orientation (Fig. 1b). We designated the direction of highest

54  *lacZ* expression as "forward". On average, "reverse" transcription neared half the "forward" activity

55  (Fig. 1c). We repeated the experiment with 25 well-characterised promoter DNA fragments[21].

56  Importantly, we selected only promoters that did not contain a detectable TSS on the opposite DNA

57  strand[19,22,23]. Such DNA fragments could only drive *lacZ* expression in the "forward" orientation

58  (Extended Data Fig. 2). For an arbitrary subset of TSS pairs, we mapped RNA 5' ends (Fig. 1d). This

59  allowed annotation of promoter elements (Fig. 1e). For the *yib*A2 DNA fragment, TSSs were

60  convergent. For all other DNA fragments, TSSs were divergent. Hence, promoter elements for

61  oppositely oriented transcripts mapped, partially or completely, to the same section of DNA (Fig. 1e

62  and Extended Data Fig. 3). Mutation of these shared promoter sequences (Fig. 1e and Extended Data

63  Fig. 3) reduced expression in both orientations (Fig. 1f).

64

65  *Bidirectional promoter sequences are widespread and obey precise organisational rules*

66  To understand global patterns of divergent transcription we analysed TSSs mapped by RNA 5'

67  polyphosphatase sequencing (PPP-seq), dRNA-seq or cappable-seq[19,22,23]. Oppositely orientated TSSs

68  tended to co-locate (Fig. 2a). To increase sensitivity, we merged the datasets (Fig. 2a, combined). This

69  identified 5,292 divergent TSSs, defined as being separated by between 25 and 7 bp; 19 % of all

70  detected TSSs in *E. coli* (Table S1). We refer to the associated promoter sequences as "bidirectional"

71  and the corresponding TSSs as "divergent TSS pairs". The most common distance between divergent

72  TSS pairs was 18 bp (Fig. 2a, top expansion). This corresponds to transcription initiation, on opposite

73 DNA strands, either side of the same promoter -10 region (Fig. 2a, top expansion). Presumably,

74 promoter element symmetry must play a major role in creating bidirectional promoter sequences. To

75 test this, we made a position weight matrix (PWM) describing all *E. coli* promoter sequences. The PWM

76 was aligned with its reverse complement across a range of spacings (i.e. altered stagger between

77 forwards and reverse PWM). We then calculated a symmetry score for each spacing. If the PWM

78 resembled the same section of DNA in both orientations, the symmetry score increased. There was

79 strong correlation between experimentally detected divergent TSS pairs and those predicted on the basis

80 of symmetry score ($R^2 = 0.85$; Fig. 2a bottom expansion). Consistent with this, a DNA sequence logo

81 generated by aligning divergent TSSs, separated by 18 bp, was symmetrical (Fig. 2b). Contrastingly,

82 TSSs with no divergent transcript generated an asymmetrical motif (Fig. 2c). Recall that the second and

83 sixth positions of $\sigma^{70}$ promoter -10 elements are crucial[5]. Non-template strand bases at these positions,

84 relative to the orientation of RNAP binding, are sequestered by $\sigma^{70}$ to stabilise initial duplex melting[5].

85 At divergent TSS pairs offset by 18 bp, these key bases reciprocally coincide on opposite DNA strands.

86 Hence, these positions are strongly conserved (Fig. 2b). An example of a -10 element with such

87 symmetry is shown in Fig. 2d; the critical bases at positions two and six are underlined. Divergent

88 transcription was also enriched for TSS pairs offset by 23, 12, 10 or 7 bp (Extended Data Fig. 4a). These

89 configurations also correspond to reciprocal base pairing between key -10 element nucleotides and

90 TSSs (Extended Data Fig. 4b). We note that, of the divergent TSS pairs detected within horizontally

91 acquired sequences, two match one of the common configurations (Fig. 1). For instance, the divergent

92 TSSs within *yigG* are 7 bp apart.  To test whether the symmetrical sequences were intrinsically able to

93 drive divergent transcription we used *in vitro* transcription assays. Bidirectional promoter sequences

94 were cloned, in either the forward or reverse orientation, upstream of the λ*oop* terminator in plasmid

95 pSR. Transcripts terminated by λ*oop* have a defined length and can be detected using electrophoresis.

96 In all cases, regardless of the cloning orientation, bidirectional promoter sequences produced detectable

97 transcripts terminated by λ*oop* (Extended Data Fig. 4c,d). Since *in vitro* transcription assays use no

98 protein factors other than RNAP, divergent transcription must be an intrinsic property of bidirectional

99 promoter DNA sequences.

100

101 *Molecular basis for promoter sequence bidirectionality: a dual role for transcription start sites*

102 In *E. coli* transcription preferentially initiates at an adenine (Fig. 2c). For divergent TSSs 18 bp apart,

103 the +1 nucleotide corresponds to position -18 on the opposite DNA strand. Hence, -18 is most often a

104 thymine (Fig. 2b). A thymine at position -18 can increase transcription by altering DNA bending[24]. This

105 change in DNA conformation enhances the interaction between the nearby $\sigma^{70}$ residue R451 and the

106 DNA backbone[24] (Fig. 3a). Importantly, this can negate the need for a -35 element[25]. We speculated

107 that the +1/-18 overlap could explain why this configuration is most frequently detected. To test this,

108 we cloned a bidirectional promoter sequence, with 18 bp between TSSs, in both orientations upstream

109    of the λ*oop* transcriptional terminator (Fig. 3bi). We also made derivatives where the A•T at each +1/-

110    18 position was replaced with C•G (Fig. 3bii-iii). Note that cloning in both orientations was necessary

111    because transcription directed away from λ*oop* is not precisely terminated. Hence, a discrete transcript

112    is not produced. As expected, altering the TSS reduced production of the associated RNA (Fig. 3c,

113    compare lane 1 with 5 and 3 with 11); the same mutations also reduced transcription in the opposite

114    direction (compare lane 1 with 9 and 3 with 7). Though σ[70] RA451 was defective at the bidirectional

115    promoter sequences (even lane numbers to 12) it was unimpaired at a control promoter not requiring

116    this contact (lanes 13-14).

117

118    *Bidirectional promoter sequences are overrepresented at sites of mRNA synthesis*

119    Of all divergent TSS pairs, 48% located to intergenic DNA (Fig. 2e). Of the resulting transcripts, 75 %

120    are expected to be mRNAs, based on the orientation of the flanking genes. By comparison, only 29 %

121    of directional TSSs were in intergenic regions, near a gene 5' end, with 89 % of associated transcripts

122    expected to be mRNAs (Extended Data Fig. 5d). This suggests many bidirectional promoter sequences

123    control gene expression. Hence, we determined if divergent TSS pairs mapped to well-characterised

124    promoters, known to control mRNA production, listed in RegulonDB[26]. First, we examined the

125    RegulonDB set of 317 mRNA TSSs identified by 5' RACE[27]. Of these TSSs, 311 were in our combined

126    TSS dataset; a 156-fold enrichment compared to random genome co-ordinates. Enrichment was

127    significantly more pronounced for divergent TSS pairs (252-fold) than directional TSSs (136-fold) (*P*

128    = 0.002, Fisher's exact test). RegulonDB lists a further 3,330 mRNA TSSs, identified using numerous

129    approaches, according to RNAP σ factor specificity[26]. The majority are σ[70] dependent[26]. Of the 1,994

130    σ[70] dependent TSSs, 1,410 are in our combined TSS dataset (a 113-fold enrichment). Moreover, σ[70]

131    dependent TSSs are significantly overrepresented amongst divergent TSS pairs (133-fold enrichment)

132    compared to directional TSSs (108-fold enrichment) (*P* = 0.032, Fisher's exact test). Conversely,

133    RegulonDB described promoters for alternative σ factors do not preferentially map to divergent TSS

134    pairs (Table S2). This is consistent with divergent TSS pairs mapping to sequences resembling the σ[70]

135    -10 element (Fig. 2b and Extended Data Fig. 4b).

136

137    *Length and stability of transcripts arising from bidirectional promoter sequences*

138    To investigate length and stability of transcripts from bidirectional promoter sequences we used RNA-

139    seq. We focused on divergent TSS pairs in non-coding regions; overlapping mRNA synthesis confounds

140    analysis of intragenic promoters. After grouping intergenic loci according to adjacent gene orientation,

141    we generated aggregate RNA coverage plots (Fig. 3d and 3e). At bidirectional promoter sequences

142    between co-oriented genes, RNAs generated in each direction had different properties (Fig. 3d). Whilst

143    non-coding antisense transcripts were detectable, coding transcripts were more abundant and longer.

144    For bidirectional promoter sequences between divergent genes, two coding RNAs are expected. Hence,

145 transcript abundance and length was similar in both directions (Fig. 3e). Fig. 3f illustrates examples of

146 RNAs derived from divergent TSS pairs. Note that cappable-seq detects only RNA 5' ends whilst RNA-

147 seq detects all RNA sequences.

148

149 *Bidirectional promoter sequences are widespread in bacteria*

150 Widespread divergent transcription from bacterial promoters has not been reported previously.

151 However, a prior study did identify a modest number of divergent TSS pairs offset by 18 bp in

152 *Pseudomonas aeruginosa*[28]. To determine the prevalence of bidirectional promoter sequences across

153 the bacterial kingdom we analysed TSS maps for proteobacteria[19,22,23,28–31], actinobacteria[32,33], and a

154 firmicute[34]. We also mapped TSSs in an additional firmicute, *Bacillus subtilis*, using cappable-seq

155 (Extended Data Fig. 5 and Table S3). Co-localised divergent TSSs were abundant in all bacteria

156 analysed (Fig. 4a). Proteobacteria and actinobacteria were most similar; divergent TSS pairs were

157 usually offset by 18 or 19 bp as in *E. coli* (Extended Data Fig. 6). Firmicutes used the same range of -10

158 element configurations illustrated in Extended Data Fig. 4 for *E. coli* but with little preference for a

159 single arrangement (Extended Data Fig. 6).

160

161 *Bidirectional promoter sequences in archaea and bacteria are analogous*

162 Archaeal transcription is closely related to that of eukaryotes; promoters have a TATA box and B

163 recognition element (BRE; 5'-CGAAA-3'), located a narrow range of distances from the TSS[35].

164 Previously, Grünberger and co-workers noted divergent transcription from sites either side of a shared

165 TATA box in *Pyrococcus furiosus*[36]. This resembles the scenario presented here for bacteria. We

166 speculated that bidirectional promoter sequences should be widespread in archaea with multiple spacing

167 preferences evident. We analysed TSS maps for the archaea *Thermococcus kodakarensis* and *Haloferax*

168 *volcanii*[37,38]. We observed strong signatures of promoter sequence bidirectionality (Fig. 4a). In *T.*

169 *kodakarensis*, divergent TSS pairs were predominantly offset by 52 bp, and located either side of a

170 shared TATA box element (5'-TTATAAA-3') (Fig. 4b,c and Extended Data Fig. 7a). Less frequently,

171 divergent TSS pairs were offset by 36 bp (Fig. 4b and Extended Data Fig. 7a). In this situation, the BRE

172 is positioned so the initial C•G bp can also act as the TSS on the opposite DNA strand (Fig. 4c). Similar

173 observations were made for *H. volcanii* despite the unusual TATA box consensus (5'-TTWT-3') of

174 haloarchaea (Extended Data Fig. 7b,c).

175

176 *Promoter sequences acquired by horizontal gene transfer are more frequently bidirectional*

177 In eukaryotes, bidirectional promoters occur more frequently in recently acquired DNA[16]. We have

178 shown that horizontally acquired bacterial genes, by virtue of their high AT-content, are enriched for

179 promoter -10 elements and TSSs[19,39]. We reasoned that many such sites could represent bidirectional

180 promoter sequences. Indeed, our initial analysis of 33 σ[70] binding events, within horizontally acquired

181 DNA, is consistent with this view (Fig. 1). As predicted, detection of divergent TSS pairs using PPP-

182 seq, increased in cells lacking H-NS; a protein that suppresses transcription at horizontally acquired

183 DNA (Extended Data Fig. 8a). Parallel DNA sequence analysis demonstrated elevated promoter

184 symmetry in foreign genes (Extended Data Fig. 8b). To understand other bacteria we utilised the TSS

185 datasets described above. For both directional and bidirectional promoter sequences we determined the

186 percentage of associated TSSs mapping to horizontally acquired sections of the cognate genome.

187 Bidirectional promoter sequences were enriched in horizontally acquired regions for 6 of the 8 genomes

188 analysed (Extended Data Fig. 8c).

189

190 *Bidirectional promoter sequences allow coordinated regulation of divergent operons*

191 The widespread occurrence of bidirectional promoter sequences has implications for our understanding

192 of gene regulation. In the bacterium *Vibrio cholerae*, the genes VC1303 and VC1304 encode a para-

193 aminobenzoate synthetase and a fumarate hydratase respectively. The divergent coding sequences share

194 the same gene regulatory region (Fig. 5a). Examination reveals a divergent transcription start site pair

195 with 23 bp spacing; the second most common configuration in both *E. coli* and *V. cholerae* (Fig. 5a and

196 Extended Data Fig. 4). Here, reciprocal base pairing is observed between -10 element positions one and

197 two (underlined in Fig. 5a). The intergenic region is also a target for the cyclic-di-GMP responsive

198 transcription factor VpsT (identified using ChIP-seq, data to be presented elsewhere). Binding of VpsT

199 was confirmed using DNAseI footprinting. As expected, in the absence of cyclic-di-GMP, VpsT was

200 unable to bind the regulatory DNA (Fig. 5b, lanes 1-5). Conversely, in the presence of cyclic-di-GMP,

201 VpsT protected a ~50 bp section of DNA from digestion (Fig. 5b, lanes 6-10). The expansion in Fig. 5a

202 illustrates that the VpsT footprint overlaps the bidirectional -10 element. To investigate the impact of

203 VpsT on transcription in each orientation, we first used *in vitro* transcription assays. We cloned the

204 regulatory DNA, in either orientation, upstream of the λ*oop* terminator in plasmid pSR. In the absence

205 of VpsT, transcripts of the expected size were detected in each orientation (Fig. 5c, lanes 1 and 3). When

206 VpsT was added, production of both transcripts was greatly reduced (Fig. 5c, lanes 2 and 4). To

207 understand the effect of VpsT *in vivo* we cloned the same promoter DNA fragment, in either orientation,

208 upstream of *lacZ* in plasmid pRW50T. In both orientations, promoter activity was significantly reduced

209 in *V. cholerae* expressing VpsT (Fig. 5d).

210

211 *RNA polymerase complexes compete at bidirectional promoter sequences*

212 At bidirectional promoter sequences RNAP can bind the same section of duplex DNA in two possible

213 orientations. This binding cannot be simultaneous; structural constraints preclude this[40]. Instead, RNAP

214 molecules likely compete to access the DNA duplex. We hypothesised that increased RNAP binding in

215 one orientation would reduce transcription in the opposite direction. Since the promoter -35 element

216 stabilises RNAP binding we introduced this sequence either side of a bidirectional -10 region. Our

217 initial attempts to clone such DNA fragments in plasmid pSR failed. Specifically, we could not isolate

218    recombinants with DNA inserts expected to generate high levels of reverse transcription. We reasoned

219    that such transcription might interfere with expression of the upstream *bla* gene. Hence, we utilised a

220    derivative of pSR with a λ*oop* terminator positioned upstream, as well as downstream, of the cloned

221    DNA. The DNA constructs generated are shown in Fig. 5e. The presence of two transcriptional

222    terminators allowed simultaneous detection of both forward and reverse RNA products following *in*

223    *vitro* transcription (Fig. 5f, lane 1) dependent on σ[70] side chain R451 (lane 2). Addition of a -35 element

224    upstream of the -10 sequence increased transcription in the forward direction (lane 3, lower band).

225    Concurrently, transcription in the reverse direction was reduced (lane 3, upper band). The inverse result

226    was obtained if the -35 element was introduced downstream of the -10 region (compare lanes 3 and 5).

227    When both -35 elements were present levels of divergent transcription increased in both directions (lane

228    7). However, increases were smaller than those detected with individual -35 elements (compare lanes

229    3, 5 and 7). Note that promoter -35 elements removed the requirement for σ[70] residue R451 (compare

230    lane 2 with lanes 4, 6 and 8). Indeed, the σ[70] R451A derivative was moderately more active in such

231    instances. Most likely, the R451-DNA contact hinders escape from near consensus promoters.

232

233    **DISCUSSION**

234    We demonstrate that divergent transcription from promoter sequences is a process conserved in all

235    domains of life. The phenomenon is similarly frequent in diverse prokaryotes (Extended Data Fig. 9)

236    and superficially resembles the situation in eukaryotes. However, the mechanistic basis is

237    fundamentally different (Fig. 6). In eukaryotes, chromosomal regions associated with divergent

238    transcription are large; bidirectionality is generated by nucleosome-depleted DNA and fortuitous

239    binding of transcriptional activators[12–15]. Hence, divergent transcripts initiate from easily

240    distinguishable sites separated by hundreds or thousands of base pairs, with no distance optimal.

241    Accordingly, each TSS is associated with a distinct RNAP binding event involving non-overlapping

242    DNA regions[9,12,13,16]. By contrast, bidirectional promoter sequences in bacteria have inherent symmetry.

243    Hence, RNAP can bind the same section of duplex DNA in either orientation. Our global TSS analysis

244    shows that symmetrical -10 elements are the main driver of divergent transcription (Fig. 2). This is

245    consistent with the unique role of this promoter motif. Thus, whilst other promoter sequences stabilise

246    RNAP binding, the -10 element also facilitates DNA opening and transcription initiation. Accordingly,

247    ancillary promoter sequences are ineffective without an appropriately positioned -10 motif. We show

248    that -10 elements, with inherent symmetry, can function independently to drive divergent transcription

249    (Fig. 3c and Extended Data Fig. 4c). In the most common situation, the +1 and -18 positions on opposite

250    strands align. This enhances the ability σ[70] side chain R451 to stabilise RNAP binding. Interestingly,

251    one example of a bidirectional -35 element was identified within the horizontally acquired *ygaQ* gene

252    (Fig. 1e). We speculate that such configurations are more likely to arise in foreign DNA; the high AT-

253    content ensures many potential -10 sequences are available. As in bacteria, divergent TSS pairs in

254    archaea are separated by preferred distances, corresponding to key bases for transcription initiation
255    overlapping on opposite DNA strands (Fig. 4c and Extended Data Fig. 7). The separation of TSSs by
256    34 bp and 36 bp in *H. volcanii* and *T. kodakarensis* respectively corresponds to alignment of the BRE
257    (important for RNAP binding) and the +1 site of initiation on the opposite DNA strand. This is similar
258    to alignment of positions -18 and +1 in bacteria.

259

260    Remarkably, despite the differences between prokaryotes and eukaryotes, our data suggest divergent
261    transcription is often a property of newly acquired DNA in both kingdoms. Thus, nascent promoters
262    can be inherently bidirectional. In bacteria, this is likely a consequence of both the DNA motif for
263    divergent transcription, and horizontally acquired loci, having a high AT-content[19]. The abundance of
264    non-canonical promoter elements is also likely to play a role[41]. Most sites of transcription within
265    horizontally acquired genes are associated with non-coding RNA production. However, bidirectional
266    promoter sequences elsewhere drive mRNA synthesis. Indeed, compared to directional promoters,
267    divergent TSS pairs are more frequently found in intergenic regions, particularly between divergent
268    genes (Fig. 2e and Extended Data Fig. 5d). Hence, divergent transcription must also have important
269    implications for gene expression. For instance, we show that transcriptional repressors can co-regulate
270    divergent operons by binding sites that overlap a bidirectional promoter sequence (Fig. 6). We also
271    show that frequency of transcription in a given orientation impacts divergent RNA synthesis (Fig. 5f).
272    Hence, bidirectional promoter sequences have inbuilt regulatory properties. Speculatively, divergent
273    transcription could also displace adjacently bound transcription factors or generate asRNAs impacting
274    adjacent genes. In conclusion, the widespread occurrence of bidirectional promoter sequences has
275    important implications for understanding gene regulation in all prokaryotes.

276

277    **MATERIALS AND METHODS**
278    *Strains, plasmids and oligonucleotides*
279    All strains plasmids and oligonucleotides used are listed in Table S4. Standard procedures for strain and
280    DNA manipulation were used throughout. All bacterial cultures were grown in LB media.

281

282    *β-galactosidase assays*
283    Assays were done according to the method of Miller[42]. Cells were grown in LB media supplemented
284    with appropriate antibiotics to mid-log phase. Values shown are the mean of three independent
285    experiments. Error bars represent the standard deviation of three independent experiments. Promoters
286    were characterised as active if they stimulated β-galactosidase activity >2-fold over background levels
287    generated by promoterless *lacZ*.

288

289    *Identification of transcription start sites by primer extension*

290 Transcript start sites were mapped for individual promoters using primer extension as described by
291 Haycocks and Grainger[43]. The RNA was purified from indicated *E. coli* strains carrying different DNA
292 fragments cloned in pRW50. The 5' end-labelled primer D49724, which anneals downstream of the
293 *Hin*dIII site in pRW50, was used in all experiments. Primer extension products were analysed on
294 denaturing 6% polyacrylamide gels, calibrated with size standards, and visualized using a Fuji phosphor
295 screen and Bio-Rad Molecular Imager FX.
296
297 *Genome-wide identification of divergent transcription start site pairs*
298 Divergent TSS pairs at bidirectional promoter sequences were identified by calculating the distance
299 between each TSS on the top and bottom DNA strands. The TSS were classified as divergent pairs if
300 the bottom strand TSS was between 7 and 25 bp upstream of the top strand TSS. If a TSS on a given
301 DNA strand could couple with multiple TSSs on the opposite DNA strand these were each counted as
302 separate TSS pairs. Similarly, if directional promoter sequences were associated with multiple TSSs
303 these also were individually counted. To compare TSSs in wild type *E. coli*, and the Δ*hns* derivative,
304 we used our previously generated data[19] and remapped TSSs. This was done using TSSpredator (version
305 1.06)[44] with the following settings: step height 0.1, step height reduction 0.09, step factor 1.5, step factor
306 reduction 0.5, enrichment factor 3, normalisation percentile 0.9, enrichment normalisation percentile
307 0.5, UTR length 300 and antisense UTR length 100. Cluster method was set to HIGHEST and all other
308 parameters were set to 0. For all other datasets, we used TSS locations provided by the original studies.
309 We designated TSSs as likely to drive mRNA synthesis if they were intergenic and in the correct
310 orientation upstream of a gene. Note that previous PPP-seq analysis[19] was done according to the
311 protocol of Singh and Wade[45].
312
313 *Promoter symmetry scoring*
314 To determine symmetry scores, we derived a PWM corresponding to sequences from -100 to +50 bp
315 relative to each TSSs for each species test. We refer to this as the "forward PWM". (Note that for the
316 heatmap in Fig. 2a, the forward PWM was derived from sequences from -100 to +100 to facilitate
317 analysis over a longer range of spacings; importantly, this does not impact the calculated scores). We
318 then made a "reverse PWM" that corresponds to the reverse complement of the forward PWM, but was
319 limited to sequences from -37 to +5 relative to the TSSs, since this is the range that includes all key
320 promoter elements for all species tested. We aligned the forward and reverse PWMs across all possible
321 spacing combinations. For each spacing, we calculated a symmetry score by (i) multiplying the fraction
322 of each of the four nucleotides at each position of the forward PWM with the fraction of each of the
323 complementary nucleotides at the overlapping position of the reverse PWM, and (ii) multiplying this
324 value by 4, taking the log (base 2), and summing for all positions within the overlapping PWM
325 positions. Stated $R^2$ values are Pearson product-moment correlation coefficients generated by
326 comparing symmetry scores with TSS spacing abundance across the spacing range shown. Symmetry

327    scores were also calculated for individual *E. coli* promoter sequences, to compare promoter sequences

328    in horizontally acquired versus non-horizontally acquired regions, and to compare promoter sequences

329    in H-NS-bound versus unbound regions. In these cases, we analysed individual promoter regions from

330    position -100 to +50 relative to the TSS. We aligned the reverse PWM for *E. coli* (derived as described

331    above) with each promoter sequence across all possible spacings. For each spacing, we determined the

332    frequency of the nucleotide found in the promoter with the corresponding nucleotide frequency in the

333    reverse PWM. We then multiplied these values for every position within the PWM. The final symmetry

334    score for each promoter sequence was calculated as the maximum score across all possible spacings

335    multiplied by a constant (to avoid extremely small numbers).

336

337    *Promoter sequence analysis*

338    To determine the distance between TSSs and promoter -10 elements we searched for the sequence 5'-

339    TANNNT-3' in the 17 bp region upstream of the TSS. If this sequence did not occur, or occurred

340    multiple times, the TSS was excluded to avoid ambiguities. To generate DNA sequence motifs we used

341    Weblogo[46]. For directional *E. coli* promoters we created two alignments, anchored by either the position

342    of the TSS or -10 element, that were then spliced together in the intervening DNA. This was required

343    because the spacing between the +1 and -10 entities is variable (Extended Data Fig. 5a) and results in

344    improper alignment unless taken into account (compare Fig. 2c and Extended Data Fig. 5b). This

345    adjustment was not required for bidirectional promoters with TSSs separated by 18 bp (Fig. 2b). In this

346    situation, juxtaposition of the TSSs and -10 elements are "locked" in place in accordance with Fig. 3

347    and the associated description.

348

349    *Proteins*

350    The *V. cholerae* RNAP holoenzyme was purified as described previously[47]. To facilitate

351    overexpression, *vpsT* was cloned in pET28a and the resulting construct used to transform T7 express

352    cells. Resulting transformants were used to inoculate 20 ml of LB media that was incubated overnight

353    with shaking at 37 °C. Subsequently, this culture was used to inoculate 2 L of fresh LB media. The

354    resulting culture was incubating at 37 °C with shaking until the $OD_{650}$ reached 0.8. Overexpression of

355    the encoded $His_6$-VpsT fusing was induced with 1 mM IPTG for 16 hours at 18 °C. Cells were then

356    recovered by centrifugation, resuspended in 25 mM Tris-HCl pH 7.5, 550 mM NaCl, 20 mM Imidazole,

357    and lysed by sonication. The cleared lysate was applied to a HisTrap (Amersham) column and bound

358    proteins eluted with an imidazole concentration gradient up to 500 mM. Fractions containing VpsT

359    were pooled and transferred into 25 mM Tris-HCl pH 7.5, 100 mM NaCl, 5 % (*v/v*) glycerol by dialysis.

360    Contaminating proteins were removed using a HiTrap heparin HP (Amersham) column and pure $His_6$-

361    VpsT was eluted with concentration gradient up to 1 M NaCl. Fractions containing the pure protein

362    were pooled and concentrated to 1 mg/ml using a Vivaspin (Sartorius) concentrator. Precipitated protein

363    was removed by filtration and the $His_6$ tag removed by thrombin digestion. The cleaved tag was

364    separated from VpsT in a final HisTrap chromatography step. The pure VpsT was concentrated to 1

365    mg/ml and glycerol added to a final concentration of 50 % (*v/v*) for storage.

366

367    *in vitro transcription assays*

368    *In vitro* transcription reactions used the method of Kolb *et al.*[48] as described by Savery *et al.*[49]. Plasmid

369    template DNA was isolated from *E. coli* transformed with pSR containing the appropriate promoter

370    DNA fragment. Reaction buffer contained 20 mM Tris pH 7.9, 200 mM GTP/ATP/CTP, 10 mM UTP,

371    5 μCi (α$^{32}$P) UTP, 500 mM DTT, 5 mM MgCl$_2$, 100 μg ml$^{-1}$ BSA and 0.2 mM cAMP. Template DNA

372    (at a final concentration of 16 μg ml$^{-1}$) was incubated with RNAP holoenzyme, derived from either *E.*

373    *coli* or *V. cholerae* as appropriate, to start the reaction.

374

375    *Comparison with mRNA transcription start sites listed by RegulonDB*

376    Matches between the combined TSS list, and TSSs listed in RegulonDB, were identified using the

377    COUNTIF function in Microsoft Excel. When matching TSSs in RegulonDB to the combined list of *E.*

378    *coli* TSSs we allowed a +/- 2 bp leeway. This was done because positions of equivalent TSSs, identified

379    using different methods, often vary slightly. Additionally, RNAP can initiate transcription from a single

380    promoter at one of several adjacent nucleotides. To calculate fold enrichment we first determined how

381    many known TSSs listed in RegulonDB matched TSSs in our combined dataset. We then determined

382    how many matches were identified if the positions of TSSs in our combined dataset were randomised

383    to any position in the genome. To calculate the stated fold enrichment the former result was divided by

384    the latter. We used the same approach with subsets of the combined TSS dataset corresponding to

385    directional or bidirectional promoter sequences. We note that, in all cases, similar results were obtained

386    if the positions of TSSs in our combined dataset were instead randomised to equivalent genomic

387    contexts (i.e. intragenic and intergenic promoters to coding and non-coding regions respectively). This

388    was expected because 19 % of TSSs in the combined list were in intragenic regions. This is comparable

389    to the 12 % of the *E. coli* genome annotated as non-coding. To test for significant enrichment of different

390    TSSs groups in RegulonDB, amongst the combined list of TSSs presented here, we used a

391    hypergeometric test. For this test, the number of successful draws was the number of RegulonDB TSSs

392    identified amongst the set of divergent 5,292 TSSs or the set of 23,813 directional TSSs (i.e. the total

393    number of draws). The population size was 4,639,675 (i.e. the number of bp in the *E. coli* U00096.2

394    genome) and the number of successes in the population was the number of TSSs in the RegulonDB

395    group being tested. To determine if there was a significant difference between the number of

396    RegulonDB TSSs found in the lists of divergent TSSs and directional TSSs we used Fisher's exact test.

397    Our null hypothesis was no difference in enrichment. To determine the expected number of RegulonDB

398    promoters amongst the bidirectional TSSs, we calculated the relative frequency of RegulonDB

399    promoters in the set of 23,813 directional TSSs. We then multiplied the relative frequency value by

400    5,292 (the number of divergent TSSs). These values were then compared to the experimental data.

402 *Mapping global transcript abundance by RNA-seq*

403 Duplicate cultures of *E. coli* strain MG1655 were grown until mid-exponential phase in LB media with

404 shaking at 37 ºC. Cells were harvested by centrifugation, flash frozen in liquid nitrogen, and lysed by

405 RNAsnap[50]. Total RNA was then purified from lysates using the Qiagen Mini RNeasy kit. Library

406 preparation and sequencing was done by Vertis Biotechnologie AG. Briefly, RNA molecules were

407 fragmented by sonication before ligation to oligonucleotide adapters at their 3' end. First-strand cDNA

408 synthesis was done using M-MLV reverse transcriptase and the 3' adapter as primer. The first-strand

409 cDNA was purified and the 5' Illumina TruSeq adapter was attached at the 3' end. After PCR

410 amplification the cDNA was purified using the Agencourt AMPure XP kit (Beckman Coulter

411 Genomics) and analysed by capillary electrophoresis. Libraries were sequenced on an Illumina Nextseq

412 500 system with a read length of 75 bp. Fastq files were deposited in Array Express (accession number

413 E-MTAB-9655). Individual sequence reads were mapped using Bowtie2[51]. The reference genome for

414 *E. coli* was that assigned Genbank accession number U00096.3. Resulting Binary Alignment Map

415 (BAM) files were used to generate wiggle plots using bam2wig.py[52,53]. These data, were used to

416 generate the aggregate plots shown in Fig. 3. For each dataset, we calculated the 10 % trimmed mean

417 of the read depth, in 10 bp bins, across all 3 kb regions centred on selected TSSs. We focused our

418 analysis on TSSs in non-coding regions to avoid confounding signals from overlapping mRNA

419 transcripts.

420

421 *Identification of transcription start sites by cappable-seq*

422 To map TSSs globally we used cappable-seq. Duplicate cultures of *B. subtilis* strain 168 ca were grown

423 until mid-exponential phase in LB media with shaking at 37 ºC. Cells were harvested and flash frozen

424 in liquid nitrogen. Total RNA was isolated as described previously with the exception that RNA

425 concentration and quality was determined on an Agilent 2200 Tapestation following the manufacturer's

426 instructions[54]. Library preparation and sequencing was done by Vertis Biotechnologie AG according to

427 the protocol described by Ettwiller *et al*.[22]. Briefly, 5' triphosphorylated RNA was capped with 3'-

428 desthiobiotin-TEG-guanosine 5' triphosphate (DTBGTP) using Vaccinia capping enzyme (New

429 England Biolabs). Biotinylated RNA was captured and eluted from streptavidin beads to obtain 5'

430 fragments of primary transcripts. These transcripts were poly(A) tailed with poly(A) polymerase before

431 conversion of the 5' CAP moiety to a 5' monophosphate using CAP-clip Acid pyrophosphatase

432 (Cellscript). An RNA adapter was ligated to the 5' monophosphate and cDNA synthesis was done with

433 an oligo(dT)-adapter primer and M-MLV reverse transcriptase. cDNAs were amplified by PCR to a

434 final concentration of 10-20 ng $\mu$l$^{-1}$. Full length cDNAs were fragmented and immobilised with

435 streptavidin magnetic beads for blunting and ligation of the 3' Illumina sequencing adapter. The

436 immobilised cDNA fragments were amplified via PCR. The sample libraries were mixed in equimolar

437 amounts 200-500 bp fragments were purified from an agarose gel after electrophoresis. The libraries

438    were sequenced on an Illumina Nextseq 500 system with a read length of 75 bp. Individual sequence

439    reads were mapped using Bowtie2[51]. The *B. subtilis* reference genome was that assigned Genbank

440    accession numbers NC_000964.3. Resulting Binary Alignment Map (BAM) files were used to generate

441    wiggle plots using bam2wig.py[52,53]. For each strand of the chromosome, we assigned TSSs to base

442    positions where the read depth increased more than 3-fold, compared to the previous base, in both

443    experimental replicates.

444

445    *DNAse I footprinting*

446    DNA fragments were excised from pSR using *AatII* and *Hin*dIII. After end-labelling using $\gamma^{32}$-ATP and

447    T4 PNK (NEB), footprints were done as previously described in buffer containing 40 mM Tris acetate

448    pH 7.9, 50 mM KCl, 5 mM $MgCl_2$, 500 µM DTT  and 12.5 µg/ml Herring Sperm DNA[47]. Resulting

449    DNA fragments were analysed on a 6 % denaturing gel. Subsequently, dried gels were exposed to a

450    Biorad phosphorscreen that was scanned using a Biorad Personal Molecular Imager.

451

452    *Assignment of transcription start sites to horizontally acquired DNA*

453    To identify horizontally acquired genomic regions in different bacteria we used DarkHorse with genus

454    level phylogenetic granularity[55]. Sections of DNA with high or low H-NS binding were identified using

455    the ChIP-seq analysis of Kahramanoglou *et al*[56].

456

457    **DATA AVAILABILITY**

458    The data that support these findings are available from the corresponding author on request. The *E. coli*

459    RNA-seq, and *B. subtilis* cappable-seq, data are available in Array Express using accession numbers E-

460    MTAB-9655 and E-MTAB-8582 respectively.

461

**FIGURE LEGENDS**

**Figure 1: Transcription start site pairs within horizontally acquired genes.** a) β-galactosidase activity derived from cryptic RNAP binding sites. Data are presented as mean values (n = 3 independent experiments) +/- SD and individual data points are overlaid as dot plots. b) Direction of transcription from cloned DNA fragments. c) Average forward or reverse β-galactosidase activity of all DNA fragments. d) Start sites mapped by primer extension for selected DNA fragments (orientations labelled a or b). Primer extension products in lanes 1 to 10, sizes in nucleotides (nt). Lanes 11-14 are sequencing reactions for calibration. e) Schematic representation of transcription start site pairs. Core promoter element sequences in the forward or reverse orientation are indicated by solid or open rectangles respectively. Speckled shading indicates converge of promoter elements on the same section of DNA. Transcription start sites shown as bent arrows. The positions of mutations (x) or deletions (Δ) are indicated. f) Effect of mutating shared core promoter elements. Data are presented as in panel a.

**Figure 2: Widespread divergent transcription from bidirectional promoter sequences in** *Escherichia coli.* a) Heatmaps made using global transcription start site (TSS) data [19,22,23] or position weight matrix analysis. TSSs on the top chromosome strand are aligned at the centre of the heatmap (bent arrow, labelled +1). Heatmap colour indicates abundance of bottom strand TSSs at that position.

The expansion shows the occurrence of bottom strand TSSs in a 50 bp window either side of all top strand promoters. b) Predominant DNA sequence motif associated with bidirectional or c) directional promoters. The x-axis break indicates the variable distance between -10 element and TSS at directional promoters. Each sequence motif was generated from 638 aligned promoters. d) a bidirectional promoter sequence between the *E. coli pfs* and *dgt* genes. TSSs are in uppercase. Promoter -10 elements are bold. Key sites of -10 element symmetry are underlined and correspond to the strongly conserved bases in panel b. The non-template strand bases at these positions, relative to the direction of transcription, are sequestered by $\sigma^{70}$ to stabilise initial DNA unwinding[5]. e) Categorisation of bidirectional *E. coli* promoters according to nearby gene organisation. Percentages indicate the proportion of bidirectional promoters in each genomic context. For comparison, 89 % of the *E. coli* genome is coding whist 6 %, 3 % and 2 % is intergenic DNA between co-directional, divergent and convergent genes respectively.

**Figure 3: Reciprocal stimulation between divergent transcription start sites.** a) Structure of RNAP bound to DNA (PDB: 6CA0)[57]. Relevant features labelled. b) DNA templates used for *in vitro* transcription. Sequences of promoter -10 elements (labelled) and TSSs (bent arrows) are shown. Plasmid vector DNA is shown by black lines and opposing DNA strands of the cloned bidirectional promoter sequence are shown by teal or grey lines. Interaction of $\sigma^{70}$ R451 and the DNA backbone is indicated by dashes. Note that only transcription towards the $\lambda oop$ terminator produces an RNA of defined length, detectable as a discrete band, following electrophoresis. Hence, to detect transcription in the opposite direction, it was necessary to invert the orientation of the cloned DNA sequence. c) Products of *in vitro* transcription (using templates in panel b) using either $\sigma^{70}$ or the R451A derivative. The RNAI transcript is derived from the replication origin of the plasmid DNA template. The transcript of interest/RNAI signal intensity is 0.16, 0.05, 0.56, 0.11, 0.12, 0.06, 0.17, 0.09, 0.06, 0.05, 0.15, 0.06, 1.24 and 1.30 for lanes 1 to 14 respectively. The control promoter has the sequence 5′-TTGGCATATGAAATTTTGAGGATTATACTACACTTA-3′. A representative example of two separate experiments is shown. d,e) Aggregate profiles of transcription detected by genome-wide RNA-seq experiments. Each plot illustrates averaged sequence read depth across all 3 kb regions centred on bidirectional promoter sequences in non-coding DNA. Shaded areas of plots indicate signals above the background level f) A 17.5 kb section of the *E. coli* genome aligned with cappable-seq and RNA-seq reads mapping to the top (teal) or bottom (grey) DNA strands. Genes are denoted by red block arrows. Transcription start sites (TSSs) are denoted by gridlines and bent back arrows. Double arrow heads indicate divergent TSS pairs at bidirectional promoter sequences.

**Figure 4: Bidirectional promoter sequences are widespread in prokaryotes**. a,b) Heatmaps indicate abundance and position of TSSs on the bottom DNA strand, relative to the nearest top strand promoter (bent arrow). Species and phylogenetic relationships are indicated to left of heatmaps. c) DNA sequence motifs derived from divergent TSSs in *T. kodakarensis*.

**Figure 5: Coordinated regulation of divergent transcription units from bidirectional promoter sequences.** a) Organisation of the region between VC1303 and VC1304 in *Vibrio cholerae*. Transcription start sites are shown by bent arrows (+1) and the region footprinted by VpsT is underlined. The bidirectional promoter -10 region is bold with key positions of symmetry underlined. Position numbers indicate distances from the downstream end of the cloned DNA fragment subsequently used. b) Pattern of DNAse I digestion with or without VpsT (2, 3, 4 or 5 μM) and cyclic-di-GMP (50 μM). The gel is calibrated with a Maxam-Gilbert GA ladder. The region protected by VpsT marked by a blue bar (triangles indicate VpsT induced DNAse I hypersensitivity). A representative example of three experiments is shown. c) Transcripts generated from the VC1303-VC1304 intergenic region by RNA polymerase *in vitro* with or without 2 μM VpsT and 50 μM cyclic-di-GMP. The transcript of interest/RNAI signal intensity is 0.11, 0.02, 0.18 and 0.05 for lanes 1 to 4 respectively. A representative example of two separate experiments is shown. d) β-galactosidase activity derived from the VC1303-VC1304 intergenic region cloned in either orientation upstream of *lacZ*. Cells were supplied with VpsT

from plasmid pAMNF. Empty plasmid was used as a control. Bars are mean values (n = 3 independent experiments) +/- SD with individual data points overlaid as dot plots. *P* was derived from a two-sided paired student's *t*-test e) DNA templates to assess competition between RNA polymerase molecules during transcription *in vitro*. Promoter -10 and -35 elements are shown by black rectangles. TSSs are indicated by bent arrows. Plasmid vector DNA shown as black lines and opposing DNA strands of cloned bidirectional promoter sequences as teal or grey lines. f) Products of *in vitro* transcription using templates in panel e. The RNAI transcript is derived from the replication origin of the plasmid DNA template. The 134 nt RNA/129 nt RNA signal intensity is 0.68, not detectable, 0.09, 0.09, 4.48, 5.96, 0.39 and 0.32 in lanes 1 to 8 respectively. A representative example of two separate experiments is shown.

**Figure 6: Promoter bidirectionality has a different basis in prokaryotes and eukaryotes.**

**REFERENCES**

462  1.    Mejía-Almonte, C. *et al.* Redefining fundamental concepts of transcription initiation in
463        bacteria. *Nat. Rev. Genet.* (2020) doi:10.1038/s41576-020-0254-8.

464  2.    Browning, D. F. & Busby, S. J. W. The regulation of bacterial transcription initiation. *Nat.*
465        *Rev. Microbiol.* **2**, 57–65 (2004).

466  3.    Haberle, V. & Stark, A. Eukaryotic core promoters and the functional basis of transcription
467        initiation. *Nat. Rev. Mol. Cell Biol.* **19**, 621–637 (2018).

468  4.    Bae, B., Feklistov, A., Lass-Napiorkowska, A., Landick, R. & Darst, S. A. Structure of a
469        bacterial RNA polymerase holoenzyme open promoter complex. *Elife* **4**, (2015).

470  5.    Feklistov, A. & Darst, S. A. Structural basis for promoter -10 element recognition by the
471        bacterial RNA polymerase σ subunit. *Cell* **147**, 1257–1269 (2011).

472  6.    Kramm, K., Engel, C. & Grohmann, D. Transcription initiation factor TBP: old friend new
473        questions. *Biochem. Soc. Trans.* **47**, 411–423 (2019).

474  7.    Butler, J. E. F. The RNA polymerase II core promoter: a key component in the regulation of
475        gene expression. *Genes Dev.* **16**, 2583–2592 (2002).

476  8.    Core, L. J., Waterfall, J. J. & Lis, J. T. Nascent RNA sequencing reveals widespread pausing
477        and divergent initiation at human promoters. *Science* **322**, 1845–1848 (2008).

478  9.    Seila, A. C. *et al.* Divergent transcription from active promoters. *Science* **322**, 1849–1851
479        (2008).

480  10.   Preker, P. *et al.* RNA exosome depletion reveals transcription upstream of active human
481        promoters. *Science* **322**, 1851–1854 (2008).

482  11.   He, Y., Vogelstein, B., Velculescu, V. E., Papadopoulos, N. & Kinzler, K. W. The antisense
483        transcriptomes of human cells. *Science* **322**, 1855–1857 (2008).

484  12.   Neil, H. *et al.* Widespread bidirectional promoters are the major source of cryptic transcripts in
485        yeast. *Nature* **457**, 1038–1042 (2009).

486  13.   Scruggs, B. S. *et al.* Bidirectional Transcription Arises from Two Distinct Hubs of
487        Transcription Factor Binding and Active Chromatin. *Mol. Cell* **58**, 1101–1112 (2015).

488  14.   Rege, M. *et al.* Chromatin Dynamics and the RNA Exosome Function in Concert to Regulate
489        Transcriptional Homeostasis. *Cell Rep.* **13**, 1610–1622 (2015).

490  15.   Wu, X. & Sharp, P. A. XDivergent transcription: A driving force for new gene origination?
491        *Cell* vol. 155 990 (2013).

492 16. Jin, Y., Eser, U., Struhl, K. & Churchman, L. S. The Ground State and Evolution of Promoter
493     Region Directionality. *Cell* **170**, 889-898.e10 (2017).

494 17. Remus T. Dame, F.-Z. M. R. and D. C. G. Chromosome organization in bacteria: mechanistic
495     insights into genome structure and function. *Nat. Rev. Genet.* (2019).

496 18. Browning, D. F. & Busby, S. J. W. Local and global regulation of transcription initiation in
497     bacteria. *Nat. Rev. Microbiol.* **14**, 638–650 (2016).

498 19. Singh, S. S. *et al.* Widespread suppression of intragenic transcription initiation by H-NS.
499     *Genes Dev.* **28**, 214–219 (2014).

500 20. Mitra, P., Ghosh, G., Hafeezunnisa, M. & Sen, R. Rho Protein: Roles and Mechanisms. *Annu.*
501     *Rev. Microbiol.* **71**, 687–709 (2017).

502 21. Keseler, I. M. *et al.* The EcoCyc database: reflecting new knowledge about *Escherichia coli*
503     K-12. *Nucleic Acids Res.* **45**, D543–D550 (2017).

504 22. Ettwiller, L., Buswell, J., Yigit, E. & Schildkraut, I. A novel enrichment strategy reveals
505     unprecedented number of novel transcription start sites at single base resolution in a model
506     prokaryote and the gut microbiome. *BMC Genomics* **17**, 199 (2016).

507 23. Thomason, M. K. *et al.* Global transcriptional start site mapping using differential RNA
508     sequencing reveals novel antisense RNAs in Escherichia coli. *J. Bacteriol.* **197**, 18–28 (2015).

509 24. Singh, S. S., Typas, A., Hengge, R. & Grainger, D. C. Escherichia coli σ 70 senses sequence
510     and conformation of the promoter spacer region. *Nucleic Acids Res.* **39**, 5109–5118 (2011).

511 25. Warman, E., Forrest, D., Wade, J. T. & Grainger, D. C. Widespread divergent transcription
512     from prokaryotic promoters. *bioRxiv* **44**, 2020.01.31.928960 (2020).

513 26. Santos-Zavaleta, A. *et al.* A unified resource for transcriptional regulation in Escherichia coli
514     K-12 incorporating high-throughput-generated binding data into RegulonDB version 10.0.
515     *BMC Biol.* **16**, (2018).

516 27. Mendoza-Vargas, A. *et al.* Genome-Wide Identification of Transcription Start Sites, Promoters
517     and Transcription Factor Binding Sites in E. coli. *PLoS One* **4**, e7526 (2009).

518 28. Gill, E. E. *et al.* High-throughput detection of RNA processing in bacteria. *BMC Genomics* **19**,
519     223 (2018).

520 29. Papenfort, K., Förstner, K. U., Cong, J. P., Sharma, C. M. & Bassler, B. L. Differential RNA-
521     seq of Vibrio cholerae identifies the VqmR small RNA as a regulator of biofilm formation.
522     *Proc. Natl. Acad. Sci. U. S. A.* **112**, E766–E775 (2015).

523 30. Kröger, C. *et al.* The primary transcriptome, small RNAs and regulation of antimicrobial
524     resistance in Acinetobacter baumannii ATCC 17978. *Nucleic Acids Res.* **46**, 9684–9698
525     (2018).

526 31. Sharma, C. M. *et al.* The primary transcriptome of the major human pathogen Helicobacter
527     pylori. *Nature* **464**, 250–255 (2010).

528 32. Cortes, T. *et al.* Genome-wide Mapping of Transcriptional Start Sites Defines an Extensive
529     Leaderless Transcriptome in Mycobacterium tuberculosis. *Cell Rep.* **5**, 1121–1131 (2013).

530 33. Jeong, Y. *et al.* The dynamic transcriptional and translational landscape of the model antibiotic
531     producer Streptomyces coelicolor A3(2). *Nat. Commun.* **7**, 11605 (2016).

532 34. Fan, B. *et al.* dRNA-Seq Reveals Genomewide TSSs and Noncoding RNAs of Plant Beneficial
533     Rhizobacterium Bacillus amyloliquefaciens FZB42. *PLoS One* **10**, e0142002 (2015).

534 35. Decker, K. B. & Hinton, D. M. Transcription regulation at the core: similarities among

535    bacterial, archaeal, and eukaryotic RNA polymerases. *Annu. Rev. Microbiol.* **67**, 113–39
536    (2013).

537    36.    Grünberger, F. *et al.* Next Generation DNA-Seq and Differential RNA-Seq Allow Re-
538    annotation of the Pyrococcus furiosus DSM 3638 Genome and Provide Insights Into Archaeal
539    Antisense Transcription. *Front. Microbiol.* **10**, 1603 (2019).

540    37.    Babski, J. *et al.* Genome-wide identification of transcriptional start sites in the haloarchaeon
541    Haloferax volcanii based on differential RNA-Seq (dRNA-Seq). *BMC Genomics* **17**, 629
542    (2016).

543    38.    Jäger, D., Förstner, K. U., Sharma, C. M., Santangelo, T. J. & Reeve, J. N. Primary
544    transcriptome map of the hyperthermophilic archaeon Thermococcus kodakarensis. *BMC*
545    *Genomics* **15**, 684 (2014).

546    39.    Lamberte, L. E. *et al.* Horizontally acquired AT-rich genes in Escherichia coli cause toxicity
547    by sequestering RNA polymerase. *Nat. Microbiol.* **2**, 16249 (2017).

548    40.    Chen, J. *et al.* Stepwise Promoter Melting by Bacterial RNA Polymerase. *Mol. Cell* **78**, 275-
549    288.e6 (2020).

550    41.    Warman, E. A., Singh, S. S., Gubieda, A. G. & Grainger, D. C. A non-canonical promoter
551    element drives spurious transcription of horizontally acquired bacterial genes. *Nucleic Acids*
552    *Res.* **48**, 4891–4901 (2020).

553    42.    Miller, J. Experiments in Molecular Genetics. (1972).

554    43.    Haycocks, J. R. J. & Grainger, D. C. Unusually situated binding sites for bacterial transcription
555    factors can have hidden functionality. *PLoS One* **11**, (2016).

556    44.    Dugar, G. *et al.* High-Resolution Transcriptome Maps Reveal Strain-Specific Regulatory
557    Features of Multiple Campylobacter jejuni Isolates. *PLoS Genet.* **9**, e1003495 (2013).

558    45.    Singh, N. & Wade, J. T. Identification of regulatory RNA in bacterial genomes by genome-
559    scale mapping of transcription start sites. *Methods Mol. Biol.* **1103**, 1–10 (2014).

560    46.    Crooks, G. E., Hon, G., Chandonia, J.-M. & Brenner, S. E. WebLogo: A Sequence Logo
561    Generator. doi:10.1101/gr.849004.

562    47.    Haycocks, J. R. J. J. *et al.* The quorum sensing transcription factor AphA directly regulates
563    natural competence in Vibrio cholerae. *PLoS Genet.* **15**, e1008362 (2019).

564    48.    Kolb, A., Kotlarz, D., Kusano, S. & Ishihama, A. Selectivity of the Escherichia coli RNA
565    polymerase Eσ38 for overlapping promoters and ability to support CRP activation. *Nucleic*
566    *Acids Res.* **23**, 819–826 (1995).

567    49.    Savery, N. J. *et al.* Transcription activation at class II CRP-dependent promoters: Identification
568    of determinants in the C-terminal domain of the RNA polymerase α subunit. *EMBO J.* **17**,
569    3439–3447 (1998).

570    50.    Stead, M. B. *et al.* RNAsnap$^{TM}$: A rapid, quantitative and inexpensive, method for isolating
571    total RNA from bacteria. *Nucleic Acids Res.* **40**, e156 (2012).

572    51.    Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**,
573    357–359 (2012).

574    52.    Afgan, E. *et al.* The Galaxy platform for accessible, reproducible and collaborative biomedical
575    analyses: 2018 update. *Nucleic Acids Res.* **46**, W537–W544 (2018).

576    53.    Wang, L., Wang, S. & Li, W. RSeQC: quality control of RNA-seq experiments.
577    *Bioinformatics* **28**, 2184–2185 (2012).

578  54.  Forrest, D., James, K., Yuzenkova, Y. & Zenkin, N. Single-peptide DNA-dependent RNA
579       polymerase homologous to multi-subunit RNA polymerase. *Nat. Commun.* **8**, 15774 (2017).

580  55.  Podell, S., Gaasterland, T. & Allen, E. E. A database of phylogenetically atypical genes in
581       archaeal and bacterial genomes, identified using the DarkHorse algorithm. *BMC*
582       *Bioinformatics* **9**, 419 (2008).

583  56.  Kahramanoglou, C. *et al.* Direct and indirect effects of H-NS and Fis on global gene
584       expression control in Escherichia coli. *Nucleic Acids Res.* **39**, 2073–2091 (2011).

585  57.  Narayanan, A. *et al.* Cryo-EM structure of Escherichia coli σ 70 RNA polymerase and
586       promoter DNA complex revealed a role of σ non-conserved region during the open complex
587       formation. *J. Biol. Chem.* **293**, 7367–7375 (2018).

**ACKNOWLEDGEMENTS**

**COMPETING INTERESTS STSTEMENT**

The authors declare that there are no competing interests.

**Figure 1**

**Figure 2**

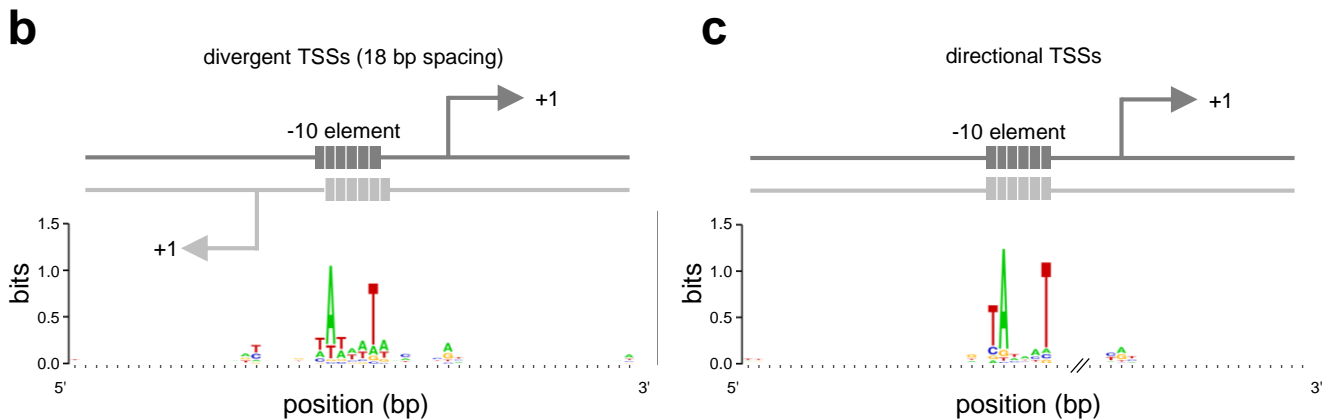**a** Position of transcription start site on top DNA strand

PPP-seq
dRNA-seq
cappable-seq
combined

-11.5 kb                                        +11.5 kb

Occurrence of nearest TSS on bottom DNA strand

maximal
medial
minimal

-50 bp          -10 element          +50 bp

combined

18 bp

prediction derived from PWM symmetry score

**b** divergent TSSs (18 bp spacing)

-10 element

bits

5'          position (bp)          3'

**c** directional TSSs

-10 element

bits

5'          position (bp)          3'

**d**

top strand RNA

read depth (PPP-seq)

pfs          dgt

bottom strand RNA

atttgaaggca**ta**g**ttt**accat**G**cgc
taaac**T**tccgta**tcaaat**ggtacgcg

**e** categorisation of 5,292 divergent TSSs according to nearby gene organisation

co-oriented (24%)          convergent (1%)

coding (51%)          divergent (24%)

# Figure 3

# Figure 4

**a**

Position of TSS on top DNA strand



**b**



**c**

**Figure 5**



**a**

atcgcagggaaagcaagag**tatcat**aaaaat**C**tcatacaaagctgctacttaaagag
tagcgtcc**C**tttcgt**tctcat**agtattttttagagtatgtttcgacgatgaatttctc

**VpsT footprint**

**b**

Promoter position (nt)

**c**

**d**

**e**

**f**

**Figure 6**



prokaryotes

eukaryotes