

Published in final edited form as:

J Mol Biol. 2007 February 16; 366(2): 420–435.

Temporal regulation of viral transcription during development of *Thermus thermophilus* bacteriophage ϕ YS40

Anastasiya Sevostyanova¹, Marko Djordjevic², Konstantin Kuznedelov³, Tatyana Naryshkina³, Mikhail S. Gelfand^{5,6}, Konstantin Severinov^{1,3,4,*}, and Leonid Minakhin^{3,*}
1From the Institute of Molecular Genetics, Russian Academy of Sciences, Moscow 123182 Russia

2From the Mathematical Biosciences Institute, The Ohio State University, Columbus, OH 43210

3From the Waksman Institute for Microbiology, Rutgers, the State University of New Jersey, Piscataway, NJ 08854

4From the Department of Molecular Biology and Biochemistry, Rutgers, the State University of New Jersey, Piscataway, NJ 08854

5From the Institute for Information Transmission Problems, Russian Academy of Sciences, Moscow 127994, Russia

6From the Faculty of Bioengineering and Bioinformatics, Lomonosov Moscow State University, Moscow, 119991, Russia

SUMMARY

Regulation of gene expression of lytic bacteriophage ϕ YS40 that infects thermophilic bacterium *Thermus thermophilus* was investigated and three temporal classes of phage genes -- early, middle, and late -- were revealed. ϕ YS40 does not encode a DNA-dependent RNA polymerase (RNAP) and must rely on host RNAP for transcription of its genes. Bioinformatic analysis using a model of *Thermus* promoters predicted 43 putative σ^A -dependent $-10/-35$ class phage promoters. A randomly chosen subset of those promoters was shown to be functional *in vivo* and *in vitro* and to belong to the early temporal class. Macroarray analysis, primer extension, and bioinformatic predictions identified 36 viral middle and late promoters. These promoters have a single common consensus element, which resembles host σ^A RNAP holoenzyme -10 promoter consensus element sequence. The mechanism responsible for the temporal control of the three classes of promoters remains unknown, since host σ^A RNAP holoenzyme-purified from either infected or uninfected cells efficiently transcribed all ϕ YS40 promoters *in vitro*. Interestingly, our data showed that during infection, there is a significant increase and decrease, respectively, of transcript amounts of host translation initiation factors IF2 and IF3. This finding, together with the fact that most middle and late ϕ YS40 transcripts were found to be leaderless, suggests that the shift to late viral gene expression may also occur at the level of mRNA translation.

*Corresponding authors Waksman Institute for Microbiology, 190 Frelinghuysen Road, Piscataway, NJ, 08854 Phones: (732) 445-6095 for K. Severinov, (732) 445-6096 for L. Minakhin FAX: (732) 445-5735 E-mail: severik@waksman.rutgers.edu, minakhin@waksman.rutgers.edu

Publisher's Disclaimer: This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final citable form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

Keywords

Thermus thermophilus; bacteriophage; bioinformatic promoter search; macroarray analysis; gene expression; leaderless mRNA

INTRODUCTION

As of the time of this writing, the complete genomic sequences of more than 380 bacteriophages (NCBI, last modified August 2006) infecting a broad variety of microorganisms have been obtained. During infection, all bacteriophages exploit resources of their hosts to redirect the host gene expression machinery to serve the needs of the virus. While comparative genomics of phages has provided important insights into the process of phage evolution, our understanding of gene expression strategies used by various phages to achieve productive infection is modest at best. Results of biochemical studies of regulation of gene expression in just a few phages (λ , T4, T7 and, more recently, XP10) have been extremely informative and provided paradigms of genetic regulation of general biological significance. For several other, less-studied, phages, recent kinetic analysis of gene transcription patterns and modeling was used to uncover viral regulatory circuits dynamics that suggested the existence of specific regulatory mechanisms^{1; 2; 3; 4}. Due to the overwhelming diversity of phages, it is clear that further studies of bacteriophage-encoded regulatory mechanisms will reveal novel paradigms of gene regulation.

Previously we presented an approach combining bioinformatics and experimental studies that allowed us to obtain a comprehensive view of temporal gene expression during infection of *Xantomonas oryzae* by phage XP10^{3; 5}. Here, we extend parts of such analysis to a much larger phage ϕ YS40 that infects hyperthermophilic eubacterium *Thermus thermophilus*. Despite recent advances in phage genomics, only a few phages infecting thermophiles have been completely sequenced to date. Most of thermophilic phages whose genomes have been determined infect hyperthermophilic archaeal species and may be of little relevance for understanding phages that infect thermophilic eubacteria^{6; 7; 8}. The subject of this study, bacteriophage ϕ YS40, is similar in its genome size⁹ and virion morphology¹⁰ to T4, a prototypical *E. coli* phage whose studies over the years revealed a staggering variety of mechanisms of regulation of gene expression. We hypothesized that like T4, ϕ YS40 may also encode a wealth of regulatory mechanisms ensuring coordinated regulation of different temporal classes of viral genes. Uncovering such mechanisms and establishing phage-encoded proteins responsible is of great interest, since proteins from thermophilic organisms are good candidates for crystallization, alone or in complex with their cellular targets. Thus, characterization of regulatory mechanisms encoded by phages infecting thermophilic bacteria will allow to approach molecular basis of genetic regulation structurally. With these ideas in mind, we studied host and viral gene expression during ϕ YS40 infection. Our results reveal temporal regulation of ϕ YS40 transcription and allow identification of early, middle and late phage promoters. Promoters from the last two temporal classes have distinct consensus elements that differ from elements of early viral and housekeeping host promoters and may define a new class of bacterial RNAP promoters. Analysis of early and middle/late phage mRNA strongly suggests that during ϕ YS40 infection there occurs a novel regulatory 'shift' from host to viral genome expression at the level of translation initiation. Thus, our results show the potential of comprehensive analysis of bacteriophage infection process for identification of novel regulatory mechanisms, and open up several new avenues for experimental investigation of genetic switches in *Thermus*.

RESULTS

Prediction of putative σ^A -dependent $-10/-35$ promoters in the ϕ YS40 genome

Bacteriophage ϕ YS40 does not encode a DNA-dependent RNA polymerase (RNAP) or any recognizable RNAP σ factor and must therefore rely entirely on host RNAP to transcribe its genes. Transcription from early ϕ YS40 promoters is most likely initiated by *T. thermophilus* RNAP holoenzyme containing the primary sigma factor, σ^A . To efficiently compete for RNAP with host promoters, early viral promoters should be strong, i.e., they are expected to have a good match to σ^A consensus promoter elements, which should allow their identification by bioinformatic means. To identify putative ϕ YS40 early promoters, we created a bioinformatic model of a *T. thermophilus* σ^A promoter. The model is based on previously reported *T. thermophilus* σ^A promoters, both those with experimentally verified transcription start points (by primer extension and/or S1 mapping) and those for which such determination was not made. Manual multiple sequence alignment of ten promoters with identified start points revealed, as expected, an unambiguous sequence conservation of the -10 and -35 promoter elements. The SignalX program¹¹ was applied to this alignment in order to make an initial positional weight matrix (profile) of *T. thermophilus* σ^A promoters. This profile assigns a numerical weight to each nucleotide at each position, so that a total score (z-score) of a candidate sequence reflects its similarity to known promoters. Five *T. thermophilus* promoters without experimentally identified start points were analyzed using the initial profile to reveal likely locations of promoter consensus elements and the final profile of a σ^A -dependent *Thermus* promoter was built using a multiple alignment of all fifteen known *T. thermophilus* promoters (see Table 1, supplementary Table S1, and Fig. 3A). The z-score of consensus *Thermus* promoter was 4.5; the highest and lowest z-scores in the training set were 4.42 and 3.02 for P₂₁₅ promoter and promoter in front of the 4.5S rRNA gene, respectively (see Table 1).

The promoter profile was used to search the ϕ YS40 genome with the GenomeExplorer program¹¹. The following search parameters were used: *i*) for every ϕ YS40 gene, a region from -200 to $+75$ bp relative to the first nucleotide of the annotated start codon was considered; *ii*) the spacer length between the -10 and the -35 promoter elements was allowed to vary between 16-18 bp; *iii*) the sequence of the spacer did not influence the search; *iv*) irrespective of its direction, a predicted promoter could intersect with an upstream gene by no more than 50 bp; and *v*) the search cutoff was set at a z-score of no less than 3.5. This cutoff was selected as a tradeoff between search specificity (absence of candidate early promoters upstream of genes coding for previously identified ϕ YS40 virion proteins⁹, which should belong to middle or late viral genes) and sensitivity (absence of likely early genes or operons without candidate early promoters in front of them).

Using these parameters, 47 putative $-10/-35$ promoter sequences were identified. For several candidate promoters, predicted transcription start points were located downstream of annotated translation start codons. Four predicted promoters for which no downstream start codons could be located were excluded from further analysis, leaving a total of 43 promoters listed in Table 2. Putative promoters for which alternative (i.e., different from those reported in the published annotation) downstream translation start codons could be located are marked by an asterisk in Table 2 (two asterisks when alternative start codons are preceded by plausible Shine-Dalgarno sequences). The new ORF coordinates are also listed in Table 2.

As expected, no promoters were predicted in non-coding regions between ϕ YS40 genes in a tail-to-tail arrangement. Among the head-to-head arranged genes (a total of 12 gene pairs), the non-coding region separating genes 27 and 28 contained two divergent predicted promoters, while the rest contained only one (gene pairs 15-16, 32-33, 36-37, 131-132, 136-137, 163-164) or no (55-56, 65-66, 95-96, 140-141, 146-147) predicted $-10/-35$ promoters. The head-to-head

transcribed regions with no predicted promoters likely contain phage promoters that are different from the $-10/-35$ class promoters (this conjecture was largely confirmed by further analysis, below).

The ϕ YS40 genome contains 170 annotated ORFs and 3 tRNA genes⁹. Two-thirds of the ϕ YS40 genome (114 genes) are transcribed in one direction (leftward in the genome map, see Fig. 1), and 56 genes are transcribed in the opposite, rightward, direction. Earlier analysis identified four gene clusters in the ϕ YS40 genome⁹. With an exception of rare “intruders”, genes within a cluster are transcribed in one direction (leftwards for cluster 1 (genes 1-36) and cluster 3 (genes 62-146), rightwards for cluster 2 (genes 37-61) and cluster 4 (genes 147-170) (Fig. 1). While clustering is statistically significant, no inferences about its functional role were made. The distribution of putative early promoters in ϕ YS40 gene clusters is highly non-random. Cluster 3 contains 30 predicted promoters, cluster 1 - 8, cluster 2 - 3, and cluster 4 - 2 promoters. In cluster 3, all putative $-10/-35$ promoters are located upstream of genes 83-137, a group of short genes that code for proteins of unknown function. In other clusters, predicted $-10/-35$ promoters are located upstream of genes involved in nucleotide metabolism, replication, recombination, and regulation of transcription⁹. Only one predicted $-10/-35$ promoter-like sequence was found upstream of a ϕ YS40 virion structural gene (gene 154), strongly indicating that a separate class of promoters is used for expression of structural (late) ϕ YS40 genes.

The logos^{12; 13} of the -35 and -10 promoter elements of *T. thermophilus* promoters and predicted ϕ YS40 early promoters are shown in Figs. 3AB. As can be seen, positions -7 , -11 , and -12 of the -10 promoter element are the most conserved ones in both the host and predicted viral promoters (the corresponding positions are also highly conserved in the *E. coli* σ^{70} -dependent promoters). Both host and viral promoters have a less conserved extended -10 “TG” motif. The -35 element of predicted phage promoters has a consensus sequence CTTGACa, compared to *T. thermophilus* cTTGACA and *E. coli* TTGACA consensus sequences. Inspection of predicted phage promoter sequences upstream of the -35 element, downstream of the -10 element, or in the spacer between the elements using the SignalX program did not reveal any additional areas of sequence similarities.

Macroarray analysis of gene expression during ϕ YS40 infection

To understand the temporal pattern of ϕ YS40 gene expression, a macroarray of ϕ YS40 genes was prepared. The array contained spots with equal amounts of PCR-amplified fragments of 29 representative viral genes. One group of spots reported the abundance of mRNA of genes from the predicted “early” region of cluster 3 (genes 83, 91, 94, 116, 131). Other spots represented genes likely involved in nucleotide metabolism, replication, and recombination (genes 15, 16, 23, 27, 32, 33, 36, 37), genes coding for structural proteins and DNA packaging enzymes (genes 1, 3, 56, 73, 82, 146, 147, 152, 163, 164), and genes coding for putative transcription regulators (genes 18 and 71). Since partially overlapping or closely spaced viral genes are likely co-transcribed (transcribed from the same promoter), some spots on the array report abundance of transcripts of multiple genes. For example, gene spots 1, 3, and 15 and 23 and 27 likely report the abundance of polycistronic mRNAs from transcription units comprising genes 1-15 and 19-27, respectively. The array also included pairs of spots corresponding to gene pairs in the “head-to-head” orientation (genes 15-16, 27-28, 32-33, 36-37, 55-56, 163-164, see Fig. 1) since these divergently transcribed genes may belong to different temporal classes (see above).

In order to determine whether ϕ YS40 shuts off host gene expression, PCR fragments corresponding to several housekeeping *T. thermophilus* genes - *rpoC* (RNAP β' subunit), *sigA* (the primary sigma factor σ^A), *dnaK* (protein chaperone), *TTHA0466* (alcohol

dehydrogenase), *infB* (translation initiation factor 2, IF2), and *infC* (translation initiation factor 3, IF3) - were included in the array. The membrane also contained spots with total genomic DNA of ϕ YS40 and its host. As a loading and normalization control, two spots containing a PCR fragment of the *zfrp8* gene from *Drosophila melanogaster* were used. *T. thermophilus* cells were infected with ϕ YS40 and total RNA was extracted 0, 25, 50, and 75 min post-infection. The time-points were selected on the basis of a single-burst experiment that indicated that a 25-min time-point corresponded to the middle of the eclipse period, the 50-min time point corresponded to its end, while at the 75-min time-point progeny phage began to be produced.

Equal amounts of total RNA from each time-point were combined with the *zfrp8* probe and used to generate radioactively labeled cDNA by random priming/reverse transcription followed by hybridization to the array. To quantitatively analyze macroarray data, radioactive signals from each spot were corrected for background and normalized based on the relative strength of the *zfrp8* spot signal. Next, the amount of radioactivity in each spot (which corresponds to transcript abundance) was plotted as a function of time post infection. As expected, the total amount of ϕ YS40 transcripts increased through infection relative to the control *zfrp8* spot (blue line in Fig. 2A). In contrast, the total amount of *T. thermophilus* transcripts normalized to the *zfrp8* spot decreased throughout the same period (red line in Fig. 2A), indicating that ϕ YS40 either shuts off host transcription or increases the rate of host transcripts decay. The abundances of some individual transcripts, such as *rpoC*, *sigA*, *infC*, and *dnaK* also decreased between the 25 and 75-min time-points (data not shown). Interestingly, the amount of the *infB* transcript, which was relatively low in the beginning of infection, increased rapidly after the 50 min time-point (a blue line in Fig. 2B), which is contrary to the rapid decrease of *infC* transcript amount during the infection (a red line in Fig. 2B). This unusual behavior is discussed in more detail in the Discussion section.

To compare the behavior of individual ϕ YS40 transcripts, plots of normalized spot signal intensity versus time post infection were scaled to make mean transcript abundances for each spot equal (Fig. 2C). Systematic clustering analysis of temporal patterns of individual ϕ YS40 genes (see Supplementary Fig. S1) revealed three different temporal classes. The averages of scaled abundances calculated for each of the three temporal classes are shown as separate panels in Fig. 2D. As can be seen, the three temporal classes are clearly distinguished by the period of time during which the greatest change in transcript abundance occurs. For the first class, significant amounts of transcripts accumulate during the first 25 min of infection. Genes from this class are classified as ϕ YS40 early genes. Transcripts of the second class have very low abundance in the first 25 min of infection but their abundance increases dramatically between 25 min and 50 min post-infection. These transcripts correspond to ϕ YS40 middle genes. Finally, the abundance of transcripts from the third temporal class is low during the first 50 min post-infection but dramatically increases afterwards. These are ϕ YS40 late transcripts.

The genomic positions of ϕ YS40 genes that belong to different temporal classes are shown in Fig. 1. Many genes with unknown function, most notably all of cluster 3 genes located downstream of predicted $-10/-35$ class promoters, belong to the early class. Genes whose products are involved in DNA replication, recombination and nucleotide metabolism also belong to this class. Every gene (or a group of likely co-transcribed genes) that behaves as early on the macroarray is preceded by a predicted $-10/-35$ class promoter, independently confirming our promoter prediction results. The only exception are co-transcribed genes 163-165. However, this group of genes is preceded by a predicted $-10/-35$ promoter with a z-score of 3.37, just below the cut-off value of 3.5 used for the search. It is therefore likely that this early promoter is functional and we therefore included it in Table 2 (marked by three asterisks).

Most ϕ YS40 middle genes encode structural proteins as well as proteins involved in DNA packaging. Late genes with known functions encode exclusively the structural proteins of the phage. There is no predicted early promoters upstream of middle and late genes revealed by macroarray analysis, again suggesting that promoters for genes of these temporal classes differ from the $-10/-35$ class promoters.

Mapping ϕ YS40 promoters *in vivo*

In our initial attempts to identify middle and late promoters of the phage, regions upstream of genes that were found to belong to the middle and late temporal classes were bioinformatically examined for the presence of common sequence motifs that were absent from the early promoters. However, no such motifs could be identified, possible due to the small number of genes examined. To experimentally identify ϕ YS40 promoters, primer extension analysis of RNA samples used in macroarray experiments was performed. Overall, 5' ends of 18 phage transcripts were identified. Primer extension product corresponding to a representative early ϕ YS40 promoter (P₈₃, Fig. 4B) peaked 25 minutes post-infection and decreased steadily afterwards. Primer extension product corresponding to a representative middle promoter P₁₄₀ appeared between 25 and 50 min post infection and steadily increased afterwards (Fig. 4B). Primer extension product corresponding to late transcripts appeared after 50 min and dramatically increased by the end of infection (a representative late transcript from P₈₂ is shown in Fig. 4B). In case of middle and late transcripts, kinetics of primer extension products accumulation during the infection matched that observed in macroarray experiments. However, primer extension products corresponding to RNA transcribed from early promoters decreased between 50 min and 75 min post-infection, while the macroarray data showed continued increase in early transcripts abundance. Since primer extension reveals abundance of mRNA transcribed from an individual promoter, it is possible that the increase of macroarray signal later in the infection is due to read-through transcription from middle and/or late promoters located further upstream. In agreement with this idea, for almost half of the early phage genes spotted on the macroarray (6 of 13) there is a predicted middle/late promoter further upstream (Fig. 1). Alternative explanations such as *i*) preferential degradation of 5' ends of early phage mRNAs or *ii*) transcription antitermination late in infection are also possible.

For each of the nine primer extension reactions designed to reveal the presence of bioinformatically-predicted ϕ YS40 $-10/-35$ promoters, expected primer extension products were obtained. Moreover, all nine promoters belonged to the early temporal class (they are underlined in Table 2). The result shows that our bioinformatic analysis identified early ϕ YS40 promoters with a high degree of confidence. Transcription start points for 6 of the early promoters were located in front of annotated genes 18, 85, 91, 103, 116, and 131 (Figs. 1 and 4A). For three other early genes, 37, 83, and 84, annotated translation start codons were located upstream of experimentally determined (and predicted) transcription start sites. However, additional start codons preceded by plausible Shine-Dalgarno motifs were found downstream of experimentally determined transcription start points, strongly indicating that initial annotations of coding sequences of these genes were incorrect. Interestingly, for 3 of the 9 $-10/-35$ promoters analyzed, transcription start points coincided with or were very close to (2 bp upstream in the case of P₁₈) translation start points (Fig. 4A).

In order to identify middle and late viral promoters, regions upstream of genes that were found to belong to the middle and late temporal classes were examined by primer extension. Nine primer extension products corresponding to seven middle (P₁₆, P₇₂, P₁₄₀, P₁₄₆, P₁₄₇, P₁₄₈, and P₁₆₄) and two late (P₅₅ and P₈₂) transcripts were identified. Sequence alignments of regions upstream of middle and late transcripts primer extension products ends revealed a common -10 -like element (consensus sequence TAA^{*}AATa) with the highest conservation of positions

-12, -11, -9, and -7 relative to the transcription start point (Fig. 3E and 4A). Also, a presence of extended -10 "TG" motif was detected in some middle/late promoters. No additional areas of conservations were apparent. Remarkably, the transcription start sites of 8 of the 9 experimentally identified middle and late promoters are located 0-2 bp upstream of the first nucleotide of annotated translation start codons (only P₁₄₈ has an obvious upstream SD motif, see Fig. 4A). Barring gross misannotation of ϕ YS40 ORF start points, the result suggests that most middle and late viral transcripts (and some early transcripts, above) are leaderless.

Since no obvious differences between the middle and late promoter sequences could be detected, a profile of a ϕ YS40 middle/late promoter was created based on an alignment of eight experimentally confirmed "leaderless" middle and late promoters. The profile included an ATG/GTG start codon located 5-10 bp downstream of the -10 element (see supplementary Table S2). The ϕ YS40 genome was searched for the presence of middle/late promoters using parameters identical to those used for early phage promoters search (see above) but the search area was limited to positions -75 to +75 relative to the first nucleotide of published annotated start codons.

37 additional candidate middle/late promoters were revealed by the search. Several predicted middle/late promoters were located inside of annotated ORFs; however, start codons associated with predicted promoters were in frame with these ORFs and could therefore be likely used for translation initiation. A few putative middle/late promoters whose ATG/GTG elements were out of frame with annotated ORFs were excluded from further analysis on the grounds that they were invariably located in areas containing predicted or experimentally confirmed early promoters (data not shown, see also below). The remaining 28 new putative middle/late promoters are listed in Table 3 (new proposed start codons that are in frame with previously annotated ORFs are marked by asterisks). Table 3 also contains eight experimentally confirmed "leaderless" middle/late promoters that were also found by the search, as expected.

To assess the quality of bioinformatic predictions of phage middle/late promoters, additional primer extension reactions were performed using primers designed to reveal predicted promoters P₅₆, P₆₅, P₈₀, and P₁₅₂. No primer extension products with P₈₀ and P₁₅₂ primers were observed. One should bear in mind, however, that the absence of primer extension products does not necessarily mean that no promoter is located in these regions, since we often find that several primers need to be tested in order to obtain a primer extension product in good yield. Most importantly, primer extension products with P₅₆ and P₆₅ primers not only matched the predicted start points but also behaved as middle transcripts (data not shown), indicating that our search reveals middle/late promoters of the phage with reasonable confidence.

Analysis of putative middle/late promoters distribution in the genome revealed the following features. First, with an exception of genes 117 and 154, genes preceded by predicted (or experimentally shown) middle/late promoters did not have predicted -10/-35 promoters in their upstream regions (in contrast, as already mentioned, putative middle/late promoters that were excluded from Table 3 on the grounds that their ATG/GTG elements were out of frame with annotated ORFs were all located in regions harboring early promoters). Second, for most divergently transcribed genes that lacked a predicted -10/-35 promoter in front of them, a putative middle/late promoter was found upstream. Third, predicted middle/late promoters were identified in front of those genes or putatively co-transcribed gene units (operons) that behaved as middle or late on the macroarray but were not tested by primer extension. Overall, the results of middle/late promoter predictions are consistent with experimental data and further extend our understanding of phage transcription. For example, consistent with the macroarray data clustering, a predicted middle-late promoter was identified in front of a rightward-

transcribed group of late genes 1-4. In the absence of such a promoter, these genes would have been grouped with early genes transcribed from the P₈-P₁₅ promoters, Fig. 1.

***In vitro* transcription from ϕ YS40 promoters**

The ϕ YS40 middle and late promoters resemble late promoters of *E. coli* bacteriophage T4 and other T4-like phages^{14; 15}. These promoters contain a single promoter element that is recognized by RNAP holoenzyme containing phage-encoded sigma factor σ^{55} ^{16; 17}. Though ϕ YS40 genome does not encode a recognizable sigma factor, it is possible that *i)* ϕ YS40-encoded σ is so divergent that it is not identified by bioinformatic means or *ii)* ϕ YS40-encoded regulators allow the σ^A RNAP holoenzyme to transcribe viral middle and late promoters at later stages of infection (or, alternatively, a phage-encoded factor prevents transcription from these promoters early in infection). To investigate this matter further and to independently confirm identification of ϕ YS40 promoters, we amplified DNA fragments containing promoters identified by primer extension *in vivo* and performed *in vitro* transcription with host RNAP σ^A holoenzymes affinity purified from ϕ YS40-infected or uninfected cells. Representative results are shown in Fig. 5. As can be seen, transcripts from both early and middle/late promoters were observed and primer extension reactions showed that in each case *in vitro* transcription start points coincided with those determined *in vivo* (data not shown). No difference in promoter utilization by RNAP purified from infected or uninfected cells was observed. Thus, the σ^A RNAP holoenzyme from uninfected cells efficiently recognized phage middle/late promoters in the absence of added factors (conversely, the σ^A RNAP holoenzyme from ϕ YS40-infected cells transcribed from early phage promoters). Likewise, *in vitro* transcription from DNA fragments containing several host promoters did not reveal any difference in transcription efficiency by RNAPs prepared from infected and uninfected *T. thermophilus* cells (data not shown).

DISCUSSION

Here, we report the results of preliminary analysis of gene expression strategy of ϕ YS40, a large bacteriophage infecting thermophilic eubacterium *Thermus thermophilus*. To our knowledge, this is the first time ever such an analysis was undertaken for any bacteriophage infecting a thermophilic bacterium. The approach that we used to identify early viral promoters involved bioinformatic analysis of phage genome for the presence of sequences with similarities to host housekeeping promoters. Primer extension and *in vitro* transcription analyses showed that our search reveals viral promoters recognized by the host σ^A RNAP holoenzyme with a high degree of confidence, and macroarray and primer extension analyses showed that these promoters belong to the early temporal class of viral genes. The predicted early phage promoters are located in front of ϕ YS40 genes that are expected (based on sequence similarities) to be expressed early in the infection. In addition, a large number of putative σ^A promoters were located in front of short genes with unknown function in the ϕ YS40 gene cluster 3. The presence of early promoters in front of these genes suggests that the products of at least some of them may be involved in host shut-off.

In general, bioinformatic predictions of bacterial promoter sequences are not highly efficient due to degeneracy of the signal. Our success in prediction of ϕ YS40 early promoters could be due to the very tight packaging of genes in the phage genomes (which increases the signal to noise ratio by limiting the length of “searchable” DNA in or close to the intergenic regions) and the fact that early phage promoters must be strong to efficiently compete with host promoters, which means that they are more similar to consensus promoters than most host promoters. Despite some differences in the content of promoter consensus elements, predicted ϕ YS40 early promoters strongly resemble host $-10/-35$ promoters, as expected. 86% of predicted phage early promoters have an optimal 17-bp spacer separating basal promoter

elements while for host promoters (both predicted and experimentally confirmed) this value is only 59%. The difference is statistically significant at least on the level of 0.1%. The optimal spacer length of most putative phage promoters may help them to efficiently compete with host promoters for the σ^A RNAP holoenzyme.

In addition to ϕ YS40 early genes, the macroarray analysis revealed the middle and late genes of the phage. By combining the information obtained by primer extension analysis of middle and late genes transcripts and by a bioinformatic search of ϕ YS40 genome, we identified ϕ YS40 middle and late promoters. Though the middle and late ϕ YS40 genes are clearly distinguished by our clustering analysis, at present we are unable to distinguish the middle and late promoters based on their sequences, and we consequently treated them together. A consensus ϕ YS40 middle/late promoter has a single promoter element that is located ~ 10 bases upstream of transcription start point and is similar but clearly distinct from the -10 consensus element of early phage (or housekeeping host) promoters.

The temporal regulation of gene expression of bacteriophage T4, a well-studied *E. coli* phage that is similar in size to ϕ YS40, is achieved by sequential interaction of host RNAP with phage-encoded proteins that change its promoter specificity (reviewed in^{14; 18}). The middle and late T4 promoters differ from early phage promoters and from each other. The middle promoters are recognized by an RNAP holoenzyme containing the primary σ factor of the host, σ^{70} , bound to phage-encoded co-activator AsiA. The middle promoters consist of an extended -10 element (consensus sequence TGnTATAAT) and an upstream MotA box to which phage-encoded co-activator MotA binds. Late T4 promoters contain a single promoter element (consensus sequence TATAAATA), which is only recognized by a holoenzyme containing phage-encoded σ factor gp55. At least *in vitro*, middle/late promoters of ϕ YS40 are efficiently recognized by *T. thermophilus* σ^A RNAP holoenzyme without any help from phage-encoded factors. This finding raises questions as to how a change in promoter specificity of host RNAP during ϕ YS40 infection is achieved. Clearly, there must exist a mechanism(s) that determines decreased utilization of early promoters late in infection and, conversely, the absence of middle/late promoter utilization early in infection. Identification of ϕ YS40 proteins that interact with host RNAP at different stages of infection may help to clarify the issue. However, *T. thermophilus* RNAP purified from ϕ YS40-infected cells using a mild single-step affinity purification procedure has unaltered promoter specificity and does not contain any proteins other than RNAP subunit based on visual inspection of Coomassie-stained gels (unpublished observations). Thus, unlike the straightforward case of T4, which encodes a number of proteins that bind host RNAP tightly, ϕ YS40 proteins that control the switch in RNAP promoter specificity may bind host RNAP weakly. Alternatively, a change in promoter specificity could be accomplished by phage-encoded DNA-binding proteins. Since the most apparent difference between host and phage early promoters and the middle/late phage promoters is the absence of the -35 consensus element in the later, it is possible that a product of an early phage gene shuts off host and early phage promoters by interacting with the -35 element and preventing its recognition by RNAP. A search for such a protein is currently ongoing in our laboratory.

Studies conducted with *E. coli* RNAP identified two classes of promoters, the $-10/-35$ class and the extended -10 class (consensus sequence TGnTATAAT). For the latter class of promoters, the properly positioned TG motif is strictly required for promoter function^{19; 20}. Since most ϕ YS40 middle/late promoters do not have such a motif, a question arises what determines their highly efficient utilization by the σ^A holoenzyme, since the -10 consensus promoter element, TATAAT, is not sufficient for promoter utilization. Recent analysis identified an additional element recognized by *Thermus* σ^A RNAP, a downstream element GGA that allows the recognition of the -10 element in the absence of either the -35 element or the TG motif²¹. However, the downstream element is absent from ϕ YS40 promoters. Closer analysis of middle/late promoters of ϕ YS40 reveals that a TG motif is present in most of them,

though its distance from the -10 element varies from 4 to 0 base pairs. SELEX experiments aimed at determining DNA sequences that strongly bind the *E. coli* σ^{70} RNAP holoenzyme revealed that fragments containing a TGTGnTATAAT sequence bind RNAP most efficiently²². On the other hand, analysis of single- and double-stranded DNAs that specifically interact with, respectively, *Thermus* σ^A and σ^A RNAP holoenzyme, indicated that a TG motif present immediately upstream of the -10 element increases the binding efficiency²¹. Thus, it is possible that TG dinucleotide located at different distances from the -10 element may make the ϕ YS40 middle/late promoters function as an extended -10 element. On the other hand, several predicted (and experimentally verified) ϕ YS40 middle/late promoters lack a TG motif. It is therefore possible that the difference in sequence of the -10 element of the middle/late promoters (consensus sequence TAaAATa) and the early promoters (consensus sequence TAtnnT) allows promoter recognition in the absence of additional basal promoter elements. Alternatively, some unrecognized sequence elements may allow the middle/late promoter function and also determine their activation at an appropriate time during infection. Mutational analysis of middle/late promoters coupled with *in vitro* transcription in the presence of extracts of infected cells collected at different times post-infection will be needed to resolve these issues.

The most striking feature revealed by our analysis of middle/late transcripts of ϕ YS40 is the fact that most of them appear to be leaderless. In fact, we were only successful in identifying middle/late promoters by including the initiating codon ATG/GTG into the search profile along with the -10 element consensus sequence. Searches using middle/late promoter profiles in the absence of a requirement for a closely located start codon tended to find phage early promoters as well as many clearly irrelevant sequences (recall that unlike its host, ϕ YS40 genome is AT rich⁹). The set of promoters revealed by our search likely includes a majority of phage middle/late promoters. However, one should bear in mind that the “leaderless” model constrain excluded middle promoters like P₁₄₈ from which mRNAs containing canonical Shine-Dalgarno sequences is transcribed (these promoters, however, are a clear minority of phage middle/late promoters).

It is formally possible that the ATG/GTG motif included in the profile of ϕ YS40 middle/late promoters functions as a basal promoter element together with the -10 promoter element. This hypothesis appears unlikely though, since biologically plausible middle/late promoters invariably contained the ATG/GTG sequence in frame with the downstream ORF. Therefore, it appears that phage middle/late transcripts are truly leaderless. In contrast, the vast majority of host as well as early phage transcripts contain Shine-Dalgarno sequences in front of their start codons and are therefore translated in a conventional way. Thus, a switch from Shine-Dalgarno-dependent to leaderless mRNA translation initiation may occur during ϕ YS40 infection.

Translation of most prokaryotic mRNAs is initiated through the 30S ribosomal subunit, which interacts with the Shine-Dalgrano sequence of the mRNA (reviewed in²³). Initiation factors IF1, IF2, and IF3 regulate the kinetics of this process. Translation of leaderless mRNAs is initiated through an alternative pathway that involves the recognition of the 5'-terminal AUG codon by 70S ribosomes (reviewed in²⁴). Increased concentrations of IF2 enhance the efficiency of “leaderless” translation while increase of IF3 concentration decreases it^{25; 26}. In this regard, it is particularly noteworthy that while abundance of most host transcripts, including the IF3 transcript, decreased during ϕ YS40 infection, the IF2 transcript behaved as a late viral gene and its abundance increased dramatically late in infection. Assuming that the change in IF2/IF3 transcripts abundance reflects the change in the amount of respective proteins, the difference may provide a mechanism for the hypothetical switch in translational initiation mechanism during ϕ YS40 infection.

The activation of IF2 transcription during ϕ YS40 infection may occur through the same mechanism as activation of middle/late transcripts. In this regard, it would be of interest to determine if there is a difference between promoters of *T. thermophilus* genes whose transcription is activated or repressed during ϕ YS40 infection.

MATERIALS AND METHODS

Prediction of ϕ YS40 promoters

The promoter recognition profiles were constructed using SignalX¹¹ implementing the formula for positional nucleotide weights from²⁷. Identification of candidate promoters in the phage genome was done using GenomeExplorer¹¹.

Bacterial strains, phage and growth conditions

The *Thermus thermophilus* HB8 strain and the ϕ YS40 phage were generously provided by Dr. Tairo Oshima, Tokyo University of Pharmacy & Life Science. The bacterium and the phage were grown in Thermus broth (TB) medium (0.6% tryptone, 0.3% yeast extract, 0.4% NaCl, 1 mM MgCl₂, 0.5 mM CaCl₂) at 65 °C with vigorous shaking. To prepare ϕ YS40 lysates, a single plaque was resuspended in 100 μ l of TB, added to 50 ml of *T. thermophilus* culture (OD₆₀₀ 0.2), and cells were allowed to grow until complete lysis occurred (usually 16-20 hours). The lysed culture was treated with 0.5 ml of chloroform and cell debris was removed by centrifugation at 10,000 g for 10 minutes. The resultant ϕ YS40 stock ($\sim 2\text{-}4 \times 10^9$ p.f.u./ml) was stored at 4 °C and used to prepare larger amounts of phage lysate by scaling up the procedure described above.

E. coli strains XL-1Blue (New England Biolabs) and BL21(DE3)(Novagen) were used for molecular cloning and protein expression.

Total DNA purification and molecular cloning

ϕ YS40 and *T. thermophilus* HB8 total DNA were purified by extraction with phenol-chloroform and subsequent precipitation with ethanol according to²⁸.

A *T. thermophilus* HB8rpoC::10H strain containing a 10-Histidine affinity tag appended to the 3' end of the *rpoC* (which encodes the RNAP β' subunit) was constructed as follows. First, a plasmid pET21tthC_{10H} expressing the *T. th. rpoC* gene with 3' terminally located 10-Histidine tag was created by recloning the corresponding PCR-modified *rpoC*-10His gene from the pET28rpoCZTth plasmid between *NdeI* and *EcoRI* sites of pET21a (Novagen) plasmid. The pET28rpoCZTth plasmid is an expression vector bearing *rpoC* and *rpoZ* genes of *T. thermophilus* HB8 and is an intermediate created during the construction of multi-gene plasmid coexpressing *T. thermophilus* RNAP core enzyme (KK, unpublished). The *T. thermophilus rpoC* gene cloned in pET28rpoCZTth was obtained through subcloning of two PCR fragments - c1tth (2381 bp) and c2tth (2231 bp) - in the pT7Blue (Novagen) blunt-end cloning vector. The c1tth and c2tth fragments were joined via a unique *AvrII* restriction site introduced in the primers used for amplification. The sequences of primers used for amplification are available from the authors upon request. The entire *T. thermophilus rpoC* gene was cut from pT7Blue and inserted into the pET28a expression vector between the *NdeI* and *EcoRI* restriction sites.

A 750 bp HB8 genomic fragment downstream of *rpoC* sequence with primers containing engineered *SalI* and *HindIII* sites. A fragment containing thermostable kanamycin resistance cassette (*kat*)²⁹ was amplified using plasmid pMKE β gal³⁰ as a template with primers containing engineered *EcoRI* and *SalI* sites. The two PCR fragments were digested with the appropriate restriction enzymes and simultaneously ligated into *EcoRI*-*HindIII*-digested pTZ19R vector, resulting in a plasmid pTZ19kat-f. The *EcoRI*-*HindIII* fragment from this

plasmid was next cloned into appropriately digested pET21thC_{10H}. The resultant plasmid, pET21thC_{10kat-f}, contains a 10His-tagged gene *rpoC* followed by *kat* cassette, which in turn is followed by a 750 bp fragment of *T. thermophilus* chromosome downstream of *rpoC*. In order to increase efficiency of subsequent transformation into *T. thermophilus*, pET21thC_{10kat-f} was transformed into and then purified from *E. coli* K12 ER2925 Dam⁻ Dcm⁻ strain (New England Biolabs), followed by digestion with *NdeI* and *HindIII*. The restriction digestion reaction was precipitated with ethanol and used for genetic transformation of *T. thermophilus* HB8 following the procedure developed by³¹. Transformants were plated onto TB plates with 1.5% agar and 30 µg/mL kanamycin. After a 48 h incubation at 65 °C, individual kanamycin-resistant colonies were picked up and grown in liquid TB containing 10 µg/mL kanamycin, followed by extraction of total genomic DNA. The presence of required insertion downstream of *rpoC* was confirmed by PCR and DNA sequencing of amplified DNA fragments. φYS40 infected the resultant *T. thermophilus* HB8rpoC::10H strains with an efficiency comparable to that of the original HB8 strain.

Plasmid pET28Tthσ^A contains the *T. thermophilus sigA* gene cloned between the *NdeI* and *EcoRI* sites of the pET28a expression vector and was a source of N-terminally hexahistidine-tagged σ^A.

Proteins

T. thermophilus RNAP containing C-terminally decahistidine-tagged β' subunit was purified as follows. Cells were grown in TB medium with 10 µg/mL kanamycin to OD₆₀₀ 0.6-0.9, harvested by centrifugation and disrupted by sonication in buffer A (10 mM Tris-HCl, pH 8.0, 500 mM NaCl, 2 mM imidazole, pH 8.0, 5% glycerol, 0.2 mg/mL PMSF, 0.4 mg/mL pepstatin). After disruption, 0.04 mg/mL DNase I was added to cell lysate followed by a 10-min incubation on ice. After centrifugation at 15,000 g for 30 min, the cleared lysate was loaded onto a chelating Hi-Trap sepharose column (Amersham) equilibrated with Ni²⁺. The column was washed with buffer A containing 40 and 80 mM imidazole and bound protein was eluted with buffer A containing 200 mM imidazole, dialyzed against buffer B (20 mM Tris-HCl, pH 8.0, 200 mM KCl, 1 mM DTT, 0.5 mM EDTA and 50% glycerol) and stored at -20 °C. The same procedure was applied for purification of RNAP from HB8rpoC::10H cells infected with φYS40.

To purify hexahistidine-tagged *T. thermophilus* σ^A, the pET28Tthσ^A plasmid was transformed in *E. coli* BL21(DE3) cells and transformants were grown in 1 L of LB medium with kanamycin at 37 °C, induced with 1 mM IPTG, harvested by centrifugation, and disrupted by sonication in buffer A. Cleared cell lysate was loaded onto chelating H-Trap sepharose column (Amersham) equilibrated with Ni²⁺, the column was washed with buffer A containing 20 mM imidazole and hexahistidine-tagged *T. thermophilus* σ^A was eluted with buffer A containing 200 mM imidazole, dialyzed against buffer B (20 mM Tris-HCl, pH 8.0, 200 mM NaCl, 2 mM DTT and 50% glycerol) and stored at -20 °C.

Primer extension

Exponentially growing *T. thermophilus* HB8rpoC::10H cells were infected with φYS40 at multiplicity of infection of 10 and harvested at various time points after infection. At the MOI of 10 used throughout the work, the efficiency of host cell infection was always greater than 95% (i.e. less than 5% of host "survivors" were detected). Total RNA was extracted with RNeasy mini kit (Qiagen) following a procedure recommended by the manufacturer. The absolute amount of total RNA extracted from 1 ml of cell culture infected at OD₆₀₀ of 0.4 was 1.5-5 µg. For primer extension reaction, 8-10 µg of total RNA were reverse-transcribed with 100 units of SuperScript III enzyme from First-Strand Synthesis kit for RT-PCR (Invitrogen) according to the manufacturer's protocol in the presence of 1 pmol ³²P end-labeled primer. The

reactions were treated with RNase H, precipitated with ethanol and dissolved in formamide-containing loading buffer. To identify primer extension products, sequencing reaction (with the fmol DNA Cycle Sequencing kit from Promega) was performed from a corresponding PCR fragment amplified from the ϕ YS40 genome using the same end-labeled primer as that used for primer extension. The reaction products were resolved on 7% sequencing gels and visualized using a PhosphorImager (Molecular Dynamics). The sequences of the primers are available from the authors upon request.

***In vitro* transcription**

Multiple-round run-off reactions contained, in 10 μ l of standard transcription buffer (40 mM Tris HCl, pH 8.0, 40 mM KCl, 10 mM MgCl₂, 3 mM β -mercaptoethanol), 20 nM of *T. thermophilus* HB8rpoC::10H RNAP core enzyme saturated with 40 nM of *T. thermophilus* σ^A and 2-4 nM of PCR fragments containing ϕ YS40 promoters. Reactions were incubated for 10 min at 65 °C, followed by the addition of ATP, CTP, and UTP (0.2 mM each), 20 μ M GTP and 3 μ Ci of [α -³²P]GTP (3000 Ci/mmol). Reactions proceeded for 7 minutes at 65 °C and were terminated by the addition of an equal volume of formamide-containing loading buffer. The reaction products were resolved on a 7% denaturing polyacrylamide gel and visualized using a PhosphorImager.

In vitro transcription reactions for subsequent primer extension analysis contained, in 50 μ l of transcription buffer, 40 nM of *T. thermophilus* RNAP core enzyme, 80 nM of *T. thermophilus* σ^A and 6-12 nM of PCR fragments containing ϕ YS40 promoters. Reactions were performed as described above, and nucleic acids were precipitated with ethanol and dissolved in RNase free water. The reaction products were then used in primer extension reactions as described above.

Macroarray membrane preparation and hybridization

DNA fragments corresponding to each of the selected ϕ YS40 genes, *T. thermophilus* HB8 housekeeping genes, and *D. melanogaster zfrp8* gene were amplified from corresponding genomic DNA using gene-specific primer pairs. The sequences of the primers are available from the authors upon request. Membrane preparation, cDNA synthesis and macroarray hybridization were performed according to².

Macroarray data analysis

After hybridization the amount of radioactivity from each spot was quantified using PhosphorImager-generated image files that were analyzed by using the ImageQuant (Molecular Dynamics) software. The background signal was subtracted from signals corresponding to every ORF spot. To allow comparison between the signals on different membranes, the background-corrected signals were normalized relative to the average of the two *D. melanogaster zfrp8* spot signals. The normalized signals were used in further data analysis.

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

Acknowledgements

The *Thermus thermophilus* HB8 strain and the ϕ YS40 phage were generously provided by Dr. Tairo Oshima (Tokyo University of Pharmacy & Life Science). We are grateful to Dr. J. Berenguer (Universidad Autonoma de Madrid) for plasmid pMKE β gal, Dr. L. Westblade (Rockefeller University) for help with construction of the *T. thermophilus* strain HB8rpoC::10H and Dr. S. Minakhina (Waksman Institute for Microbiology) for *D. melanogaster zfrp8* gene DNA fragment. This work was supported by grants GM59295 and GM64530 from NIH (to KS). MD acknowledges support

from NSF under Agreement No. 112050 and NSF grant MCB- 8690418891. M.S.G. was partially supported by grants from the Howard Hughes Medical Institute (55005610), INTAS (05-1000008-8028), and the Russian Academy of Sciences (Program “Molecular and Cellular Biology”).

References

- Ventura M, Foley S, Bruttin A, Chennoufi SC, Canchaya C, Brussow H. Transcription mapping as a tool in phage genomics: the case of the temperate *Streptococcus thermophilus* phage Sfi21. *Virology* 2002;296:62–76. [PubMed: 12036318]
- Minakhin L, Semenova E, Liu J, Vasilov A, Severinova E, Gabisonia T, Inman R, Mushegian A, Severinov K. Genome sequence and gene expression of *Bacillus anthracis* bacteriophage Fah. *J Mol Biol* 2005;354:1–15. [PubMed: 16226766]
- Semenova E, Djordjevic M, Shraiman B, Severinov K. The tale of two RNA polymerases: transcription profiling and gene expression strategy of bacteriophage Xp10. *Mol Microbiol* 2005;55:764–77. [PubMed: 15661002]
- Duplessis M, Russell WM, Romero DA, Moineau S. Global gene expression analysis of two *Streptococcus thermophilus* bacteriophages using DNA microarray. *Virology* 2005;340:192–208. [PubMed: 16043205]
- Djordjevic M, Semenova E, Shraiman B, Severinov K. Quantitative analysis of a virulent bacteriophage transcription strategy. *Virology*. 2006In press
- Peng X, Blum H, She Q, Mallok S, Brugger K, Garrett RA, Zillig W, Prangishvili D. Sequences and replication of genomes of the archaeal rudiviruses SIRV1 and SIRV2: relationships to the archaeal lipothrixvirus SIFV and some eukaryal viruses. *Virology* 2001;291:226–34. [PubMed: 11878892]
- Haring M, Vestergaard G, Rachel R, Chen L, Garrett RA, Prangishvili D. *Virology*: independent virus development outside a host. *Nature* 2005;436:1101–2. [PubMed: 16121167]
- Prangishvili D, Garrett RA, Koonin EV. Evolutionary genomics of archaeal viruses: unique viral genomes in the third domain of life. *Virus Res* 2006;117:52–67. [PubMed: 16503363]
- Naryshkina T, Liu J, Florens L, Swanson SK, Pavlov AR, Pavlova NV, Inman R, Kozyavkin SA, Washburn M, Mushegian A, Severinov K. *Thermus thermophilus* bacteriophage φYS40 genome and proteomic characterization of virions. *J Mol Biol*. 2006In press
- Sakaki Y, Oshima T. Isolation and characterization of a bacteriophage infectious to an extreme thermophile, *Thermus thermophilus* HB8. *J Virol* 1975;15:1449–53. [PubMed: 1142476]
- Mironov AA, Vinokurova NP, Gelfand MS. Software for analysis of bacterial genomes. *Mol. Biol. (Mosk)* 2000;34:222–31.
- Schneider TD, Stephens RM. Sequence logos: a new way to display consensus sequences. *Nucleic Acids Res* 1990;18:6097–100. [PubMed: 2172928]
- Crooks GE, Hon G, Chandonia JM, Brenner SE. WebLogo: a sequence logo generator. *Genome Res* 2004;14:1188–90. [PubMed: 15173120]
- Miller ES, Kutter E, Mosig G, Arisaka F, Kunisawa T, Ruger W. Bacteriophage T4 genome. *Microbiol Mol Biol Rev* 2003;67:86–156. [PubMed: 12626685]
- Nolan JM, Petrov V, Bertrand C, Krisch HM, Karam JD. Genetic diversity among five T4-like bacteriophages. *Virol J* 2006;3:30. [PubMed: 16716236]
- Kassavetis GA, Geiduschek EP. Defining a bacteriophage T4 late promoter: bacteriophage T4 gene 55 protein suffices for directing late promoter recognition. *Proc Natl Acad Sci U S A* 1984;81:5101–5. [PubMed: 6382259]
- Kolesky SE, Ouhammouch M, Geiduschek EP. The mechanism of transcriptional activation by the topologically DNA-linked sliding clamp of bacteriophage T4. *J Mol Biol* 2002;321:767–84. [PubMed: 12206760]
- Brody EN, Kassavetis GA, Ouhammouch M, Sanders GM, Tinker RL, Geiduschek EP. Old phage, new insights: two recently recognized mechanisms of transcriptional regulation in bacteriophage T4 development. *FEMS Microbiol Lett* 1995;128:1–8. [PubMed: 7744235]
- Camacho A, Salas M. Effect of mutations in the “extended –10” motif of three *Bacillus subtilis* sigmaA-RNA polymerase-dependent promoters. *J Mol Biol* 1999;286:683–93. [PubMed: 10024443]

20. Burr T, Mitchell J, Kolb A, Minchin S, Busby S. DNA sequence elements located immediately upstream of the -10 hexamer in *Escherichia coli* promoters: a systematic study. *Nucleic Acids Res* 2000;28:1864–70. [PubMed: 10756184]
21. Feklistov A, Barinova N, Sevostyanova A, Heyduk E, Bass I, Vvedenskaya I, Kuznedelov K, Merkiene E, Stavrovskaya E, Klimasauskas S, Nikiforov V, Heyduk T, Severinov K, Kulbachinskiy A. A basal promoter element recognized by free RNA polymerase sigma subunit determines promoter recognition by RNA polymerase holoenzyme. *Mol Cell* 2006;23:97–107. [PubMed: 16798040]
22. Gaal T, Ross W, Estrem ST, Nguyen LH, Burgess RR, Gourse RL. Promoter recognition and discrimination by EsigmaS RNA polymerase. *Mol Microbiol* 2001;42:939–54. [PubMed: 11737638]
23. Laursen BS, Sorensen HP, Mortensen KK, Sperling-Petersen HU. Initiation of protein synthesis in bacteria. *Microbiol Mol Biol Rev* 2005;69:101–23. [PubMed: 15755955]
24. Moll I, Grill S, Gualerzi CO, Blasi U. Leaderless mRNAs in bacteria: surprises in ribosomal recruitment and translational control. *Mol Microbiol* 2002;43:239–46. [PubMed: 11849551]
25. Tedin K, Moll I, Grill S, Resch A, Graschopf A, Gualerzi CO, Blasi U. Translation initiation factor 3 antagonizes authentic start codon selection on leaderless mRNAs. *Mol Microbiol* 1999;31:67–77. [PubMed: 9987111]
26. Grill S, Moll I, Hasenohr D, Gualerzi CO, Blasi U. Modulation of ribosomal recruitment to 5'-terminal start codons by translation initiation factors IF2 and IF3. *FEBS Lett* 2001;495:167–71. [PubMed: 11334885]
27. Mironov AA, Koonin EV, Roytberg MA, Gelfand MS. Computer analysis of transcription regulatory patterns in completely sequenced bacterial genomes. *Nucleic Acids Res* 1999;27:2981–9. [PubMed: 10390542]
28. Sambrook, J.; Fritsch, EF.; Maniatis, T. *Molecular Cloning: A Laboratory Manual*. 2nd edit.. Cold Spring Harbor Laboratory Press; Cold Spring Harbor, NY: 1989.
29. Matsumura M, Katakura Y, Imanaka T, Aiba S. Enzymatic and nucleotide sequence studies of a kanamycin-inactivating enzyme encoded by a plasmid from thermophilic bacilli in comparison with that encoded by plasmid pUB110. *J Bacteriol* 1984;160:413–20. [PubMed: 6090428]
30. Moreno R, Zafra O, Cava F, Berenguer J. Development of a gene expression vector for *Thermus thermophilus* based on the promoter of the respiratory nitrate reductase. *Plasmid* 2003;49:2–8. [PubMed: 12583995]
31. Koyama Y, Hoshino T, Tomizuka N, Furukawa K. Genetic transformation of the extreme thermophile *Thermus thermophilus* and of other *Thermus spp.* *J Bacteriol* 1986;166:338–40. [PubMed: 3957870]
32. Hartmann RK, Erdmann VA. *Thermus thermophilus* 16S rRNA is transcribed from an isolated transcription unit. *J Bacteriol* 1989;171:2933–41. [PubMed: 2722737]

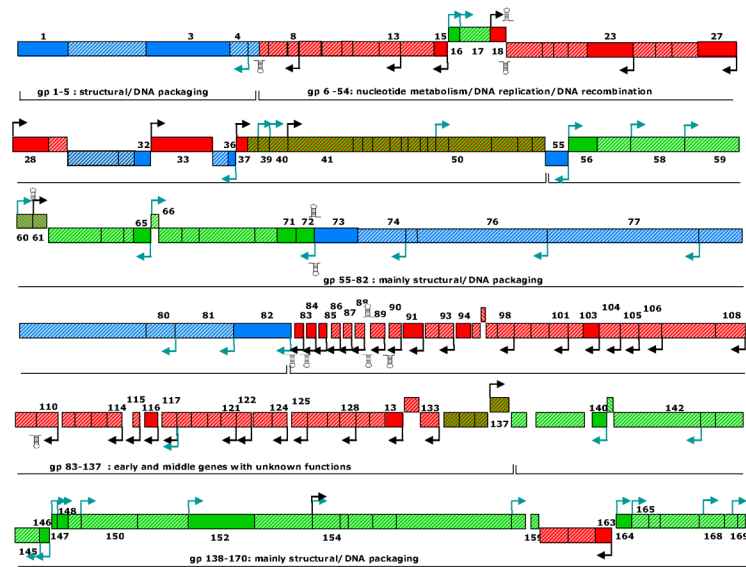


Fig 1.

Transcription map of the *T. thermophilus* bacteriophage ϕ YS40.

Colored boxes on the corresponding strand of the phage DNA represent each gene: upper boxes indicate genes with rightward orientation; lower boxes indicate genes with leftward orientation. The genes belonging to different temporal classes (defined by macroarray analysis and primer extension) are shown in different colors: early, red; middle, green; late, blue. The genes that likely belong to the corresponding classes are represented by shaded boxes of the corresponding color. Double-colored shaded boxes indicate genes with uncertain temporal class. The genes with numbers shown were used in macroarray and/or primer extension analysis or have predicted promoters. The functional modules are indicated by brackets at the bottom of the map. Promoter locations are depicted as bent arrows colored in black or blue to indicate early or middle/late promoters, respectively. Hairpins indicate possible rho-independent terminators.

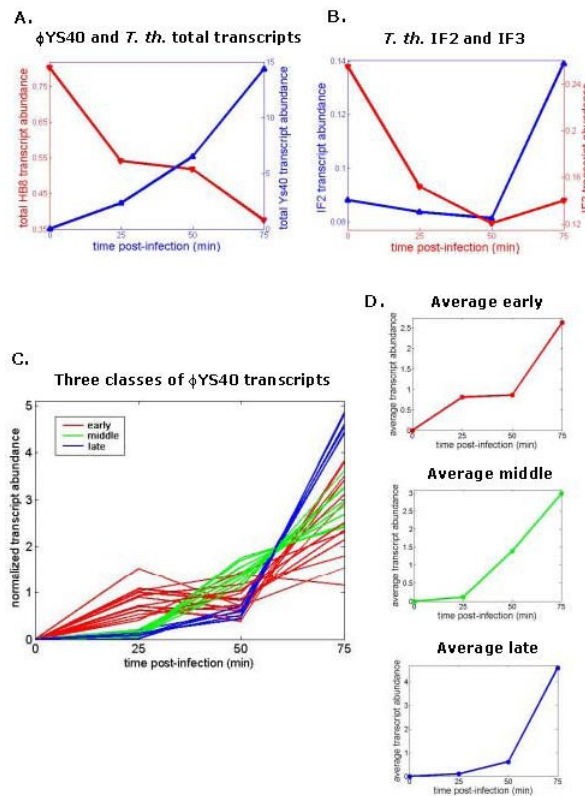


Fig 2.
Macroarray data analysis.

- A.** The abundance of total ϕ YS40-encoded transcripts (blue line) is shown together with the abundance of total *T. thermophilus*-encoded transcripts (red line).
- B.** The transcript abundances of the translation initiation factors IF2 and IF3 (blue line and red line, respectively) are shown together.
- C.** Normalized transcript abundances are presented for individual ϕ YS40 transcripts as a function of time. Transcripts that belong to different temporal classes are shown in different colors. The curves are colored according to Fig. 1: early, in red; middle, in green; late, in blue. Classification of individual transcripts into the three temporal classes is performed by the procedure described in Supplementary Appendix 1.
- D.** The three vertical panels on the right show averaged normalized transcript abundances corresponding to the three temporal classes.

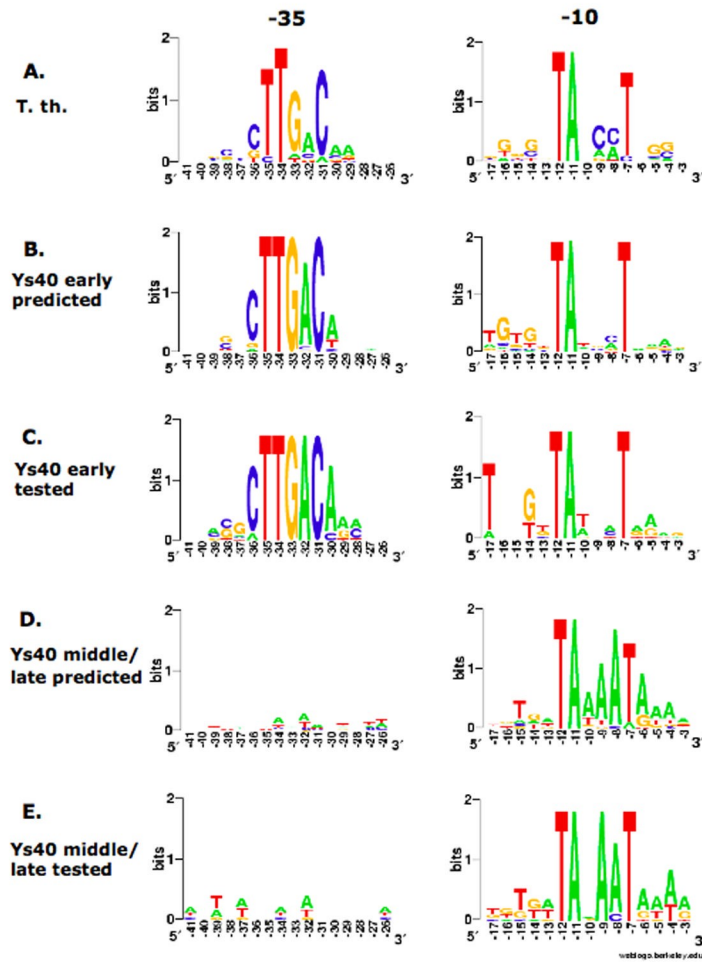


Fig. 3

Fig 3. Sequence logo representation of *T. thermophilus* and ϕ YS40 promoters. Consensus sequences were plotted with WebLogo¹³. Height of letter indicates degree of conservation. Positions are done with respect to putative or identified transcription start sites.

A. *T. thermophilus* -10/-35 promoter sequence logo with independently aligned the -35 and the -10 regions.

B-E. ϕ YS40 predicted early (B), verified early (C), predicted middle/late (D) and verified middle/late (E) independently aligned promoter consensus sequences are plotted.

A. ϕ Ys40 tested promoters

	-35	TG	-10	+1	mRNA
early promoters					
P ₁₈	TGTTTGCTCTCTAGGGGGC TTGACA ATTCTTGTCAATT TGTTTATGCT AGTCAT GCT <u>atg</u>				LL
P ₃₇	TTATTTTCCCGGTGCGTAT TTGAC AGCTATGCTGTTTT TGTTAATAT TGAACA AGG GAGAA CA atg				
P ₈₃	GGTGATTAAGTGTGGGGC TTGACA AAAAGGATTAAGTAGGG TATTAT GAAGGT G GAGGTGAGTGGG atg				
P ₈₄	ACGCTATTGGCAGATAGGC TTGACA AAATAGTCATTA TGGGTTATGAT GAAGGT G GAGGTGAGTAGG atg				
P ₈₅	TTTTCTTTTACTCCCCC TTGACA AAATGTACTAAAAGAGG TATCTT AGTGTT G GGAGGTGAAC CA atg				
P ₉₁	CATCTCAAATATATAACGC TTGACA ACGCTTAGCGCT TTGCTAACCT AAAA CA atg				LL
P ₁₀₃	CTGAGCGCTTGCACGCT ACTTTGACCT AAGCGTTTTGTCC TTAAACT AAGGC G AGGGGGTGAAGAG atg				
P ₁₁₆	ATTATTAGTGTACCAACGC TTGAC AGCGTTTCGCTAT TGTGT TAGCT TGAAG CA atg				LL
P ₁₃₁	ATTAATAGTGTACCAACGC TTGAC AGCCTTAGGTT TGT TCTATGATTTCTTCT GT § 101bp § gtg				
middle/late promoters					
P ₁₆	CACAAC TTTGGGGT GTCTGTCAACAGGTAGCAGAATTGG TGATA AAAT AAT CT G CT atg				LL
P ₅₅	TGACTTTATATATATCACAAAATTTTCGTTCAATTAAGAC TGTA AACTAATAT T AT atg				LL
P ₇₂	AAAAATGCCTCCCATTTGTTGGGGGCATTTCTTTTT TGTTATA ATGTAGCT AG CA atg				LL
P ₈₂	AAACGCTTTTTCAAATTTAAATTAATGCTAAAACATTC TGAT TAA AA TAGTACT T atg				LL
P ₁₄₀	TAAATTTCCATACACACATCATAACCACAAGACGACT TTGTTA TAGAATATAAT T atg				LL
P ₁₄₆	TTATATTGTAAAACAAAAGCGAAAAACAAAAGTTCC TGTA TAA AA TGAAAT T atg				LL
P ₁₄₇	TTTCATTTTATACAGGAAC TTTTGTTTT TCGCTTT TGTTT TACAATATAAA T gtg				LL
P ₁₄₈	ACCTATTTTCGAAAATTTTAAAGAACAAGGACTTT TGTG ATAGAATTAAGT G AGGAAGT atg				
P ₁₆₄	ACCACACCGTGCCTTTGGTGTCAAGCGCTGATGGAAGAT TGTTAA AA T AAAGCT T atg				LL

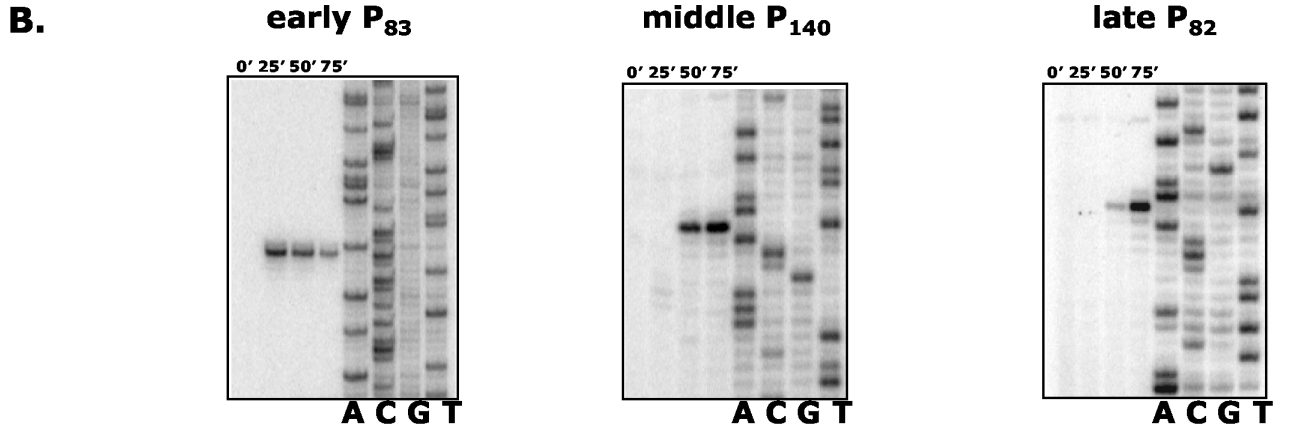


Fig 4.
 ϕ YS40 verified promoters.

- A.** Alignment of the sequences of verified ϕ YS40 promoters is shown. The -35, -10 and “TG” putative promoter elements are shown in bold. Experimentally determined transcription start sites are both boldface and underlined. The assigned translation initiation codons are shown in bold small case. Putative leaderless mRNAs transcribed from the corresponding promoters are indicated as LL.
- B.** The kinetics of accumulation of representative *in vivo* primer extension products obtained with early, middle and late phage transcripts during infection.

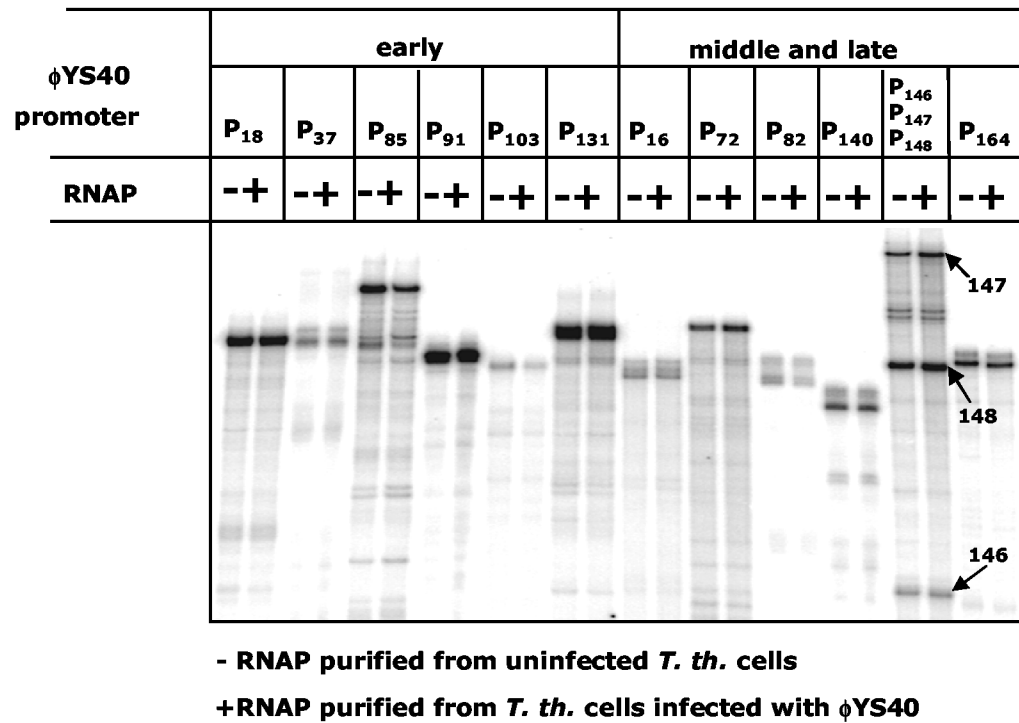


Fig 5.
Transcription by *T. thermophilus* RNAP- σ^A holoenzyme.
The results of multi-round run-off transcription from representative phage promoters by RNAP- σ^A holoenzymes purified from cells infected with the ϕ YS40 or uninfected are shown.

	Promoter sequence ^a		Accession number or reference	Spacer length	Score ^b
	-35	-10			
P31*	CGGCCCAAGC <u>TTACCA</u> AAATCCCGCCCGCTCC	<u>TAGCCT</u> GGGGCAAGSAGG	D43662	18	3.86
P35*	CCTCCTTCTG <u>TGAC</u> CGGACGGGAGGAGGGCC	<u>TATCCT</u> GGGTAAAGCTTGG	D43663	17	4.23
P37*	GGCCCTGGCC <u>TGACC</u> ATCCTCCTTGGCCT	<u>TATCCT</u> TAGGGTGGCTCCG	D43664	17	4.35
P39*	CGCCTCGCC <u>TGAC</u> CGGGAGGAGGCAACGGGG	<u>TAAAC</u> AGGGCGAGAGCG	D43665	17	3.20
P43*	GAATTCGGCT <u>TGTC</u> AGTAGCTTAGCTATGG	<u>TAACT</u> AGACCTGGAGGT	D43666	17	3.78
P214*	TAGGAGGGCC <u>TGCC</u> ATCCGCCCTTAGAGTG	<u>TACCAT</u> AGCGATTGCCCA	D43672	17	4.01
P215*	TGAGTAGCACT <u>TGAC</u> ATCATAAAGTGCTTGAGG	<u>TATCAT</u> CCGACCTGGGCGC	D43673	17	4.42
<i>Tth slp*</i>	GCCCGACCGC <u>TGACA</u> AGGGCGGTGAGGTTTT	<u>TAAGAT</u> AGCGCCGATGGG	X57333	17	3.94
<i>Tth dnaK</i>	CTCAACTCCCT <u>TGACA</u> AAATGGGCATGTGGT	<u>TAGCCT</u> GGGAGCGAGGTG	Y07826	18	4.32
<i>Tth rpsT</i>	ACCACCAAGT <u>TGCCC</u> CTTAGGGCAAGGTGTGCT	<u>TACACT</u> GGGCCCTGGTTG	AJ295159	19	3.70
<i>Tth argF*</i>	GCCTAGGCCCT <u>TGCAT</u> AAGTTGCGGGCACGGGG	<u>TATGCT</u> TAAAGCCTCATGG	Y18353	18	3.13
<i>Tth ORF4 (arg) *</i>	TTTACACCACT <u>TGACA</u> GCCTTTGATTCTGAGTC	<u>TATCCT</u> CTATTGGGGAGC	Y18353	17	4.50
<i>Tth 4.5S rRNA</i>	CTCAAGCCCT <u>TAGCCT</u> TCAGGGCTTCATGGGTG	<u>TATACT</u> ACCAGAGCCCGG	X12643	19	3.02
<i>Tth 16S rRNA</i>	TCGCAAGCCT <u>TGACA</u> AAAGGAGGGGATTGA	<u>TAGCAT</u> GGCTTTTCTGCG	Ref. 32	17	4.24
<i>Tth 23S rRNA</i>	TGGGGGCCCT <u>TGACA</u> AAAGGCATGCCTCCTTGG	<u>TATCTT</u> CCCTTTTGGCT	Ref. 32	18	4.03

* Genes with experimentally mapped transcription start sites. The sites are underlined.

^a Shown in grey boxes: -35 and -10 sequences.

^b Scores are computed using the final profile.

Table 1.
Thermus thermophilus promoters.

Table 2

Predicted early promoters of the bacteriophage ϕ YS40.

	Strand	Location ^a	Sequence and spacer ^b	Distance ^c	Score	Gene function
ORF8 **	<=	9399..9920	TTGACA-17-TATgCT	-11→10	4.20	dUTPase
ORF13	<=	12792..13367	TTGACA-18-TAaIICT	13	3.64	Recombination protein
ORF15	<=	14743..15036	TTGACA-18-TAagCT	17	3.96	Unknown
ORF18	=>	16640..17050	TTGACA-17-TATgCT	3	4.20	DNA binding HTH-domain protein, transcription regulator
ORF23	<=	19944..21620	TTGACA-17-TAgCaT	0	4.24	DNA primase bacterial DnaG type
ORF27	<=	23898..25247	TTGACt-17-TAcgAT	0	3.57	DEAD domain helicase
ORF28	=>	25396..26796	TTGACA-18-TATagT	0	3.67	Unknown
ORF33	=>	30387..32498	TTGACA-17-TAcCaT	1	4.24	DNA polymerase, without N-terminal 5'-3' exonuclease domain
ORF37 **	=>	33417..33746	TTGACA-17-TAataT	-27→10	3.56	Unknown
ORF41	=>	35201..37594	TTGACt-17-TATCaT	1	4.05	Ribonucleotide reductase, alpha subunit, N-terminal portion
ORF61 *	=>	52062..52484	TTGACt-17-TATgaT	-23→4	3.75	Unknown
ORF83 **	<=	84867..85073	TTGACA-17-TATtaT	-8→13	3.80	Unknown
ORF84 ***	<=	85328..85534	TTGACA-17-TATgaT	-11→13	4.12	Unknown
ORF85 **	<=	85767..85919	TTGACA-17-TATCT	12	4.13	Unknown
ORF86 ***	<=	86022..86258	TTGACg-17-TAagaT	-8→7	3.61	Unknown
ORF87	<=	86382..86618	TTGACt-17-TAagaT	32	3.51	Unknown
ORF88	<=	86909..87154	TTGACA-17-TATagT	30	3.67	Unknown
ORF89 **	<=	87505..87990	TTGcCA-18-TAaCCT	75	4.03	Unknown
ORF90 **	<=	88074..88439	TTGACA-17-TaggaT	-56→34	3.94	Unknown
ORF91	<=	88642..89250	TTGACA-17-TAaCCT	0	4.26	Unknown
ORF93	<=	89796..90221	TTGcCt-17-TAgCCT	12	3.72	Unknown
ORF98	<=	91835..92380	TTGACA-17-TAcaCT	13	4.08	Unknown
ORF101	<=	93619..94131	TTGACt-17-TAGCCT	11	3.95	Unknown
ORF103	<=	94885..95373	TTGACc-17-TAaaCT	13	3.87	Unknown
ORF104	<=	95510..96025	TTGACt-17-TAagCT	18	3.59	Unknown
ORF105	<=	96096..96626	TTGACc-17-TAagCT	45	3.81	Unknown
ORF106	<=	96833..97354	TTGACc-17-TAaCCT	17	4.11	Unknown
ORF108	<=	99280..100323	TTGACt-17-TAagCT	29	3.59	ATPase
ORF110	<=	101227..101973	TTGACt-17-TAggCT	35	3.65	Glycosyltransferase
ORF114	<=	103616..104107	TTGACt-17-TATCCT	12	4.13	Unknown
ORF115	<=	104451..104693	TTGACA-17-TAaggaT	14	3.94	Unknown
ORF116	<=	104803..105279	TTGACA-17-TAgcCT	0	3.95	Unknown
ORF117 *	<=	105422..105928	TTGACA-17-TATaCT	-51→0	4.26	Unknown
ORF121	<=	107552..108046	TTGACA-17-TATaaT	18	4.18	Unknown
ORF122	<=	108141..108644	TTGACA-17-TATCCT	16	4.50	Unknown
ORF124	<=	109328..109819	TTGACA-17-TATagT	21	3.67	Unknown
ORF125	<=	109998..110513	TTGACc-18-TATaaT	16	3.65	Unknown
ORF128	<=	111663..112133	TTGACA-17-TAgCaT	0	4.24	Unknown
ORF131	<=	113202..113630	TTGACA-17-TATgAT	104	4.12	Unknown
ORF133	<=	114385..115032	TTGcCA-17-TATaCT	11	3.86	Unknown
ORF137	=>	116815..117474	TTGACA-17-TATagT	65	3.67	Unknown
ORF154 ***	=>	136388..137287	TTGACA-18-TATaT	5	3.89	Unknown
ORF163	<=	146390..147022	TTGACA-17-TAaggT	6	3.37	Unknown

ORF's preceded by experimentally verified promoters are underlined.

^aLocation^a: genomic coordinates, re-annotated are shown in bold.

^b Capitals: consensus nucleotides.

^c „Distance”: the distance between the start of transcription and the start codon of the gene.

* Change of the distance by selection of the candidate start codon downstream of the annotated start of ORF.

** Change of the distance by selection of the candidate start codon with a strong Shine-Dalgarno box downstream of the annotated start of ORF.

*** Z-score of this promoter is below 3.5.

Table 3

Predicted middle/late promoters of the bacteriophage ϕ YS40.

	Strand	Location ^a	Sequence and spacer ^b	Distance ^c	Score	Gene function
ORF4	<=>	7412..8068	TAAAAA-(6)-gTG	1	3.92	unknown
ORF16	>=>	15124..15453	TAAAAA-(8)-ATG	3	4.31	unknown
ORF17	>=>	15467..16576	TAAAAA-(9)-ATG	4	4.31	IMP dehydrogenase/GMP reductase
ORF36*	<=>	33031..33318	TAAAAA-(7)-gTG	11->2	3.92	unknown
ORF39	>=>	34188..34616	TAAAAA-(5)-ATG	0	4.31	unknown
ORF40	>=>	34631..35155	TAAcATA-(7)-ATG	2	3.64	unknown
ORF50	>=>	40538..41013	TAAAATg-(5)-ATG	0	4.10	unknown
ORF52	>=>	42536..43408	TAAATA-(5)-ATG	0	3.64	N-acyltransferase
ORF55	<=>	44426..45127	TAAAcTA-(7)-ATG	2	3.92	unknown
ORF56	>=>	45187..46209	TAAATA-(9)-ATG	4	3.64	unknown
ORF58	>=>	47564..49414	TAAATt-(5)-ATG	0	3.69	unknown
ORF59	>=>	49453..51312	TAAATA-(8)-ATG	3	4.02	serine kinase
ORF60	>=>	51410..51997	TAAgATA-(8)-ATG	3	3.64	dNMP kinase
ORF65*	<=>	55466..56062	TAAATA-(7)-ATG	47->2	4.02	terminal protein in replication
ORF66	>=>	56049..56315	TAAAATg-(5)-gTG	0	3.71	unknown
ORF72	<=>	61167..61682	TAAATg-(8)-ATG	3	3.81	unknown
ORF74*	<=>	63204..64826	TAAATg-(9)-ATG	-8->4	3.81	unknown
ORF76	<=>	65085..69662	TAAAAA-(10)-ATG	5	4.31	unknown
ORF77	<=>	69684..74918	TAgAATA-(5)-ATG	0	4.13	unknown
ORF80	<=>	79880..80743	TAAAAA-(5)-gTG	0	3.92	unknown
ORF81	<=>	80788..82740	TAAATA-(9)-ATG	4	4.02	unknown
ORF82	<=>	82771..84609	TAAAAA-(6)-ATG	0	4.31	unknown
ORF117	<=>	105422..105979	TAAAAaa-(7)-ATG	2	3.64	unknown
ORF140	<=>	120226..120777	TAgAATA-(5)-ATG	0	4.13	unknown
ORF142	<=>	120953..123997	TAAAAaa-(7)-ATG	2	3.64	unknown
ORF145	<=>	125598..126548	TAAATA-(5)-ATG	0	3.64	unknown
ORF146	<=>	126553..126813	TAAAATg-(5)-ATG	0	4.10	unknown
ORF147	>=>	126870..127055	TAcAATA-(5)-gTG	0	3.63	unknown
ORF150	>=>	127979..129967	TAAAAA-(7)-ATG	2	4.31	baseplate assembly protein
ORF152	>=>	131870..134260	TAAAAaa-(6)-ATG	1	3.64	wac fibrin neck whisker
ORF154	>=>	136388..137287	TAAAATg-(5)-ATG	0	4.10	unknown
ORF159	>=>	143322..143846	TAgAATA-(8)-ATG	3	4.13	unknown
ORF164	>=>	147094..147639	TAAAAA-(5)-ATG	0	4.31	unknown
ORF165	>=>	147677..148306	TAAAATg-(8)-ATG	3	4.10	unknown
ORF168	>=>	150256..151341	TAAAATg-(5)-ATG	0	4.10	unknown
ORF169*	>=>	151284..151907	TAAATA-(5)-ATG	-54->0	4.02	unknown

ORFs preceded by experimentally verified promoters are underlined.

^a„Location“: genomic coordinates.

^b Capitals: consensus nucleotides.

^c "Distance": the distance between the start of transcription and the start codon of the gene.

* Change of the distance by selection of the candidate start codon downstream of the annotated start codon.