

In the format provided by the authors and unedited.

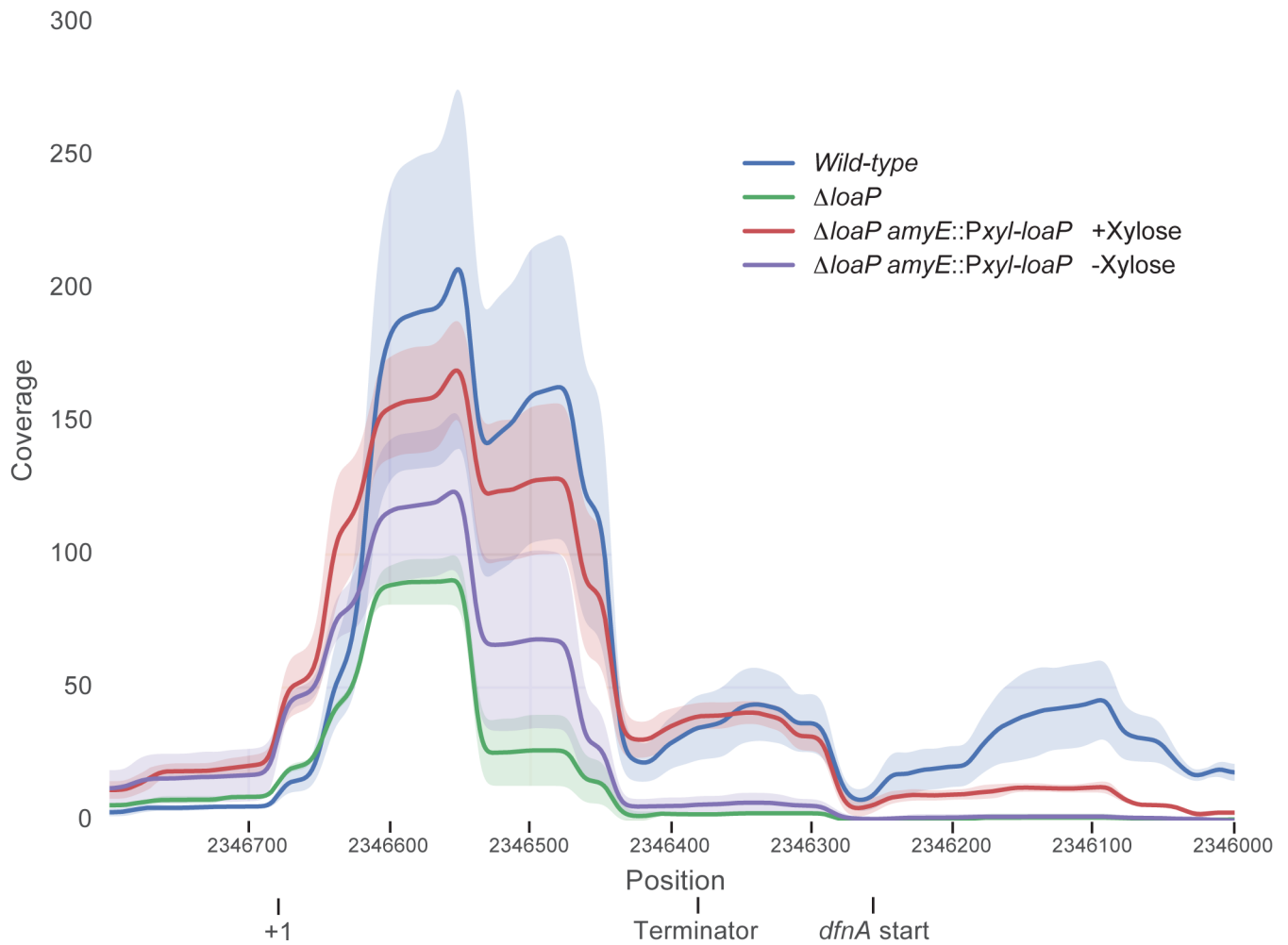
LoaP is a broadly conserved antiterminator protein that regulates antibiotic gene clusters in *Bacillus amyloliquefaciens*

Jonathan R. Goodson¹, Steven Klupt¹, Chengxi Zhang², Paul Straight^{2,3}, Wade C. Winkler^{1,3}

¹The University of Maryland, Department of Cell Biology and Molecular Genetics, College Park, MD;

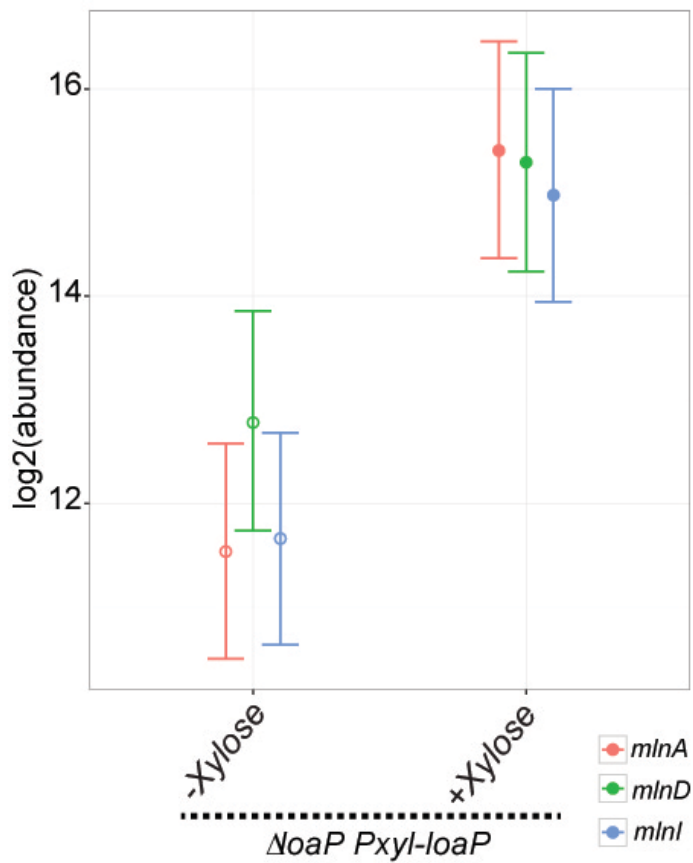
²Texas A&M University, Department of Biochemistry and Biophysics, College Station, TX

³Corresponding authors

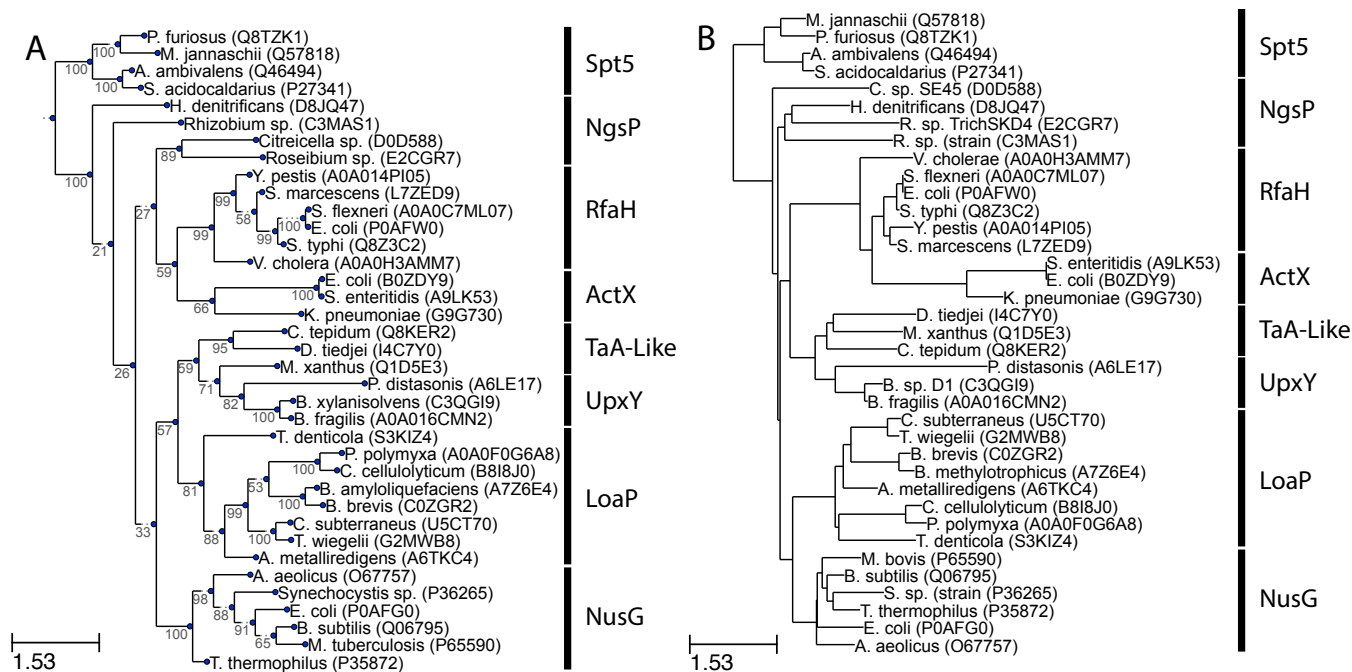


Supplementary Figure 1. Coverage of the *dfnA* leader region in RNA-seq data. RNA-seq coverage across the *dfnA* leader region normalized with DESeq2 normalization factors. Traces represent coverage data smoothed with Gaussian smoothing with a bandwidth of 5 nucleotides. Shading represents standard deviation from libraries from three independent cultures for each condition.

yfp reporter contained a constitutive promoter upstream of the *dfnA* leader region, which was transcriptionally fused to a downstream *yfp* gene. Reporter constructs were created with or without a mutation converted the UUCG tetraloop sequence to UUCA, a mutation that is predicted to abolish proper RNA hairpin formation. These data demonstrate that *B. amyloliquefaciens* LoaP antitermination can be recapitulated in the heterologous *B. subtilis* host. Data and error bars represent the mean fluorescence and 95% confidence interval (CI over means of each replicate) for all cells in three fields of view for each of three biological replicate cultures of each strain with and without induction. (E) Representative images of induced reporter strains quantified in (D) showing lack of reporter expression when the tetraloop sequence is mutated.



Supplementary Figure 3. RNA transcript levels for genes in the *mln* gene cluster are increased upon *LoaP* induction. Normalized transcript abundance at the beginning, middle, and end of the *mln* operon (*mlnA*, *mlnD*, and *mlnI*) as measured by qRT-PCR. Filled points represent samples with *loaP* expression and empty points represent samples with no or minimal *loaP* expression. Error bars represent Bayesian 95% highest posterior density estimates of mean expression. Data resulted from four biological replicate cultures.

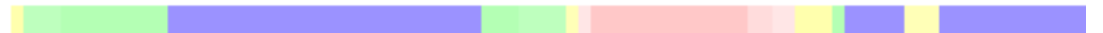


Supplementary Figure 4. Large-scale phylogenetic analysis recapitulates targeted analysis of representative sequences. (A) Phylogenetic tree constructed from multiple sequence alignment of 37 representative protein sequences. (B) Subtree consisting of the same 37 protein sequences extracted from large-scale 1,207 sequence phylogenetic analysis. Both trees show very similar topology.

EcNusG	15	25
EcRfaH	10	20
BaLoaP	9	19

E._coli_NusG	NSEA-----PKKRWYVQAFSGFEGRVATSLRE-----HIKL--HNMEDLFG--
B._subtilis_NusG	ME-----KNWYVVHTYSGYENKVKANLEK-----RVES--MGMDQKIF--
M._tuberculosis_NusG	MTTFDGDTSAGEAVDLTEANAFQDAAPAAEEVDPAAALKAELRSKPGDNWYVHVSYAGYENKVKANLET-----RVQN--LDVGDIYIF--
Synechocystis_sp._NusG	MSFTDDQSPVAEQ-----NKKTPEGHWFVAVQVASGCEKRVKLNLEQ-----RIHT--LDVADRIL--
T._thermophilus_NusG	MS-----IENYAVHTLVGQEEKAKANLEK-----RIKA--FGLQDKIF--
E._coli_RfaH	M-----QSWYLLYCKRGQLQRAQEHLER-----QAV-----
Y._pestis_RfaH	M-----KSWYLLYCKRGQLRAKEHLER-----QTV-----
S._typhi_RfaH	M-----QSWYLLYCKRGQLQRAQEHLER-----QAV-----
S._flexneri_RfaH	M-----QSWYLLYCKRGQLQRAQEHLER-----QAV-----
S._marcescens_RfaH	M-----ESWYLLYCKRGQLLRAQEHLER-----QQV-----
V._cholera_RfaH	M-----KRWYLLYCKRGEQQRKMHLEN-----QSV-----
T._denticola_LoaP	-----MDYVVVQVSTGKEKNFIEDAEFKNFDELSY-----
B._amyloliquefaciens_LoaP	-----MKWYALFVESGKEEIVQKFLRLQF--DEQAL-----
P._polymyxa_LoaP	-----VSWYVFFVVRTGREEQVKQLINEML--DSEVY-----
B._brevis_LoaP	-----LKWYVIFVESGKEEYVQKYLRLYF--NEQSL-----
C._cellulolyticum_LoaP	-----MYWYVLFVRTGREENVKKLLSKRL--DKDLF-----
T._wiegelsii_LoaP	M-----KKWYVIFTRSGYENKVKIENCFL--KQEEV-----
C._subterraneus_LoaP	M-----KKWYVLFTKSGCEEKVGKIIKKI--WENEI-----
A._metalliredigens_LoaP	-----MHVKSNEEMKAKKLVEKEI--EDI-----
B._fragilis_UpxY	-----MKSMLAAVYRLYHEKKTRDRLTA--MGI-----
B._xylanisolvenens_UpxY	MIKKNDORDLLSTDV-----IGSSVARSKRMLVAIVRICHEKKTSERLTK-----MGI-----
P._distasonis_UpxY	M-----MKKMYALKVFYRNFSEIKATLSR-----DGI-----
M._xanthus-TaA	NPGP-----RCAENDWVALLVRYNHEKYAAAQLGK-----HGY-----
D._tiedjei-TaA	NTRDVLLSRYPEE-----RPLDEDLGSWVMHCKPNCCKIASYFLS-----RNI-----
C._tepidum-TaA	NTNAL-----KKDGCWYAVYVRSRYEKKVHQYLLE-----KGL-----
Citricella_sp._NgsP	NSVLGSNGKNTSGEAVRP---YLGLRVGDAVPVEGNSVAIFDRGEIANYALLCRPOOERHAESNLAA-----RGV-----
H._denitrificans_NgsP	M-----TWYAIRTNPOREFLAGRYDE-----NGEWRPGVLEKKGYD
Roseibium_sp._NgsP	MSKMDVRYQNRALANDYD--F-----LRALVRYVSAGGMVLACCHPTKEQHALRQLTE-----RGL-----
Rhizobium_sp._NgsP	-----MGHWYVVRTRAGQQQKATREFED-----NGV-----
P._furiosus_Spt5	MA-----GKIFAVRVTHGQEETTAKLIYS-----KVR-----
M._jannaschii_Spt5	-----MIFAVRTMVGQEKNIAGLNAS-----RAE-----
S._acidocaldarii_Spt5	NEDF-----KYRNYYVLRVTGGQEIINVALILEE-----RIK-----
A._ambivalens_Spt5	MES-----KIRNYYAVKVTGGQEVSVGLMLEE-----RAK-----

cons



EcNusG Sec. Structure	
EcRfaH Sec. Structure	

44	59	82	103	117
30	45	68	88	102
33	48	71	99	113

```

EVMVPTTEEVVEI-RG--GQRR--KSERKFFPGYVLQVMVMDA-----SNHLVRSVPRVMGFIG--GT-----SDRFAPISDKEVDAINN-----RLQQV--
RVVVPPEEEETDI-KN--GKKK--VVKKKVFPGYVLVEIVMTDD-----SMYVVRNTPGVTGFVG--SAG-----SGSKPTPLLPGEAETLK-----RMGMD--
QVEVPTEEVTEI-KN--GQRK--QVNRKVLPGYILVRMDLTD-----SMAAVRNTPGVTGFVG--AT-----SRFSALALDDVVKFLL-----PRGSTRK
QVEIPKTPIVKI-RK--DGARY--QGEEKIFPGYVLIRMINDDD-----AWQVVKNTPHVINFGV--SEQKRHYGRGRGHVLFMPLSHGEVERIFR-----HYDEQ--
QVLIPTEEVVEL-RE--GGKKE--VYRKKLFPGYLFIQMDLGDDEEPEAWEVVRGTPGITGFVG--AG-----MRFVPLSPDEVRHILE-----VSGLL--

NCLAPMITLEKI-VR--GKRT--AVSEPLFPNYLFVEFDPEVIH-----TTTINATRGVSHFVR--FG-----ASF AIVPSAVIHQLSV-----YKPK--
NCWTPIVAIEKI-VR--GKRI--EVIEALFPNYLFAEFDPENIH-----TTTYSATRGGSHFVR--FG-----TQF AVIPATVIADMQA-----HAYD--
SCLTPMITLEKM-VR--GKRT--FVSEPLFPNYLFVEFDPEVIH-----TTTINATRGVSHFVR--FG-----AHP AIVPSSVIHQLSI-----YKPE--
NCLAPMITLEKI-VR--GKRT--AVSEPLFPNYLFVEFDPEVIH-----TTTISATRGVSHFVR--FG-----ASF AIVPSAVIHQLSV-----YKPK--
NCLSPITILEKI-VR--GKRI--AVSEPLFPNYLFVEFDPERIH-----TTTISATRGVSHFVR--FG-----TLF SVIPSKVIDELRT-----HASE--
ECFYPEVCVEKI-LR--GKRQ--MVQEPLFSPYMFVRFDFENGP-----SFTTVRSTRGVVDFVR--LG-----PHFRELQGLIYQLKQ-----LDCEQ--

SIVFPQRILKIR-K--AGKYT--EKQLPVFAGYLFIGTDEISKDLQY----HLRCKGQFYRFLP--NN-----QEPKFLERGFDEILNQ-----FISFG--
YSIIPKKKVTET-K--AGIKY--EALKKMFPGYVLFKTKMTERTFH-----KIKELPISCRIVN--NGAYYS--KERKTYFTTITDEEILPIIR-----LIGEG--
KFFIPLQERLFK-V--AGIVK--KEMAPLFPSPYVFIESNLPDLOFVSTNSMIC TSSDIIRLLR--YS-----KFE ASMRDSEKQML ES-----LCNDS--
KSI VPKRLVPEK-K--SGIVY--NVLKNFPGYVLIQENTNEMFH-----KIKKIPRSLRLVN--NGSYYS--QDEGAYYSSIEEKEITPILO-----LMNGG--
LPFVPLNERIFK-K--AGTVN--KNMEILFPGYVFIESKVASQEFVKMTNELMKSLODIIRLVR--YS-----DIE IAVRESERTILOS-----LYNNN--
KLLIPKRKIIER-V--KGQPY--EKIKLFPGYVVFYNAEMSDDLYY-----KISEVLKRGIFLK--EG-----KRP AFVKEEEMK IILA-----LTKNS--
EVLIPRRKIIER-I--KGEER--EKIKLFPGYVVFYKTEMTEAKYH-----EITSVLKQGVFLK--ED-----KMPASVKEEEMRVILN-----LTGDS--
KVIVPQRIPEK-R--QGETR--HWKXILFTGYLFLNVELDVTYY----KLKRIPSIHRFLG--LE-----KFEAIPLEEMQRYLR-----LCNQG--

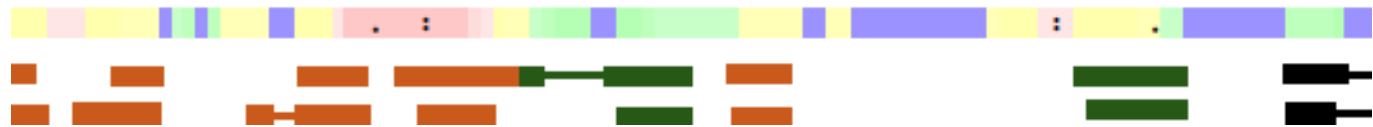
ESFLPVQEEIHO-W--SDRRK--KIERYVPMIFVHVDPAERA EVL-----TLSSVSRY--MVL R--GO-----STF AVIPDEQMERFRF-----MLDYSE--
ENFLPIQEEVHO-W--SDRRK--VVDRLVLPMMIFVHVDPOEQEVL-----TLSAISRY--MVL R--GE-----STF AVVPDQCMLRFKF-----MLDYS D--
ESYIPMKTRSYP-QP--DGEVY--IRRVLVAMLMFLRCODDYISSLN-----SILDEKAMYHRP--GT-----QIF ASIPDEEMDMFIM-----LTSSMED--

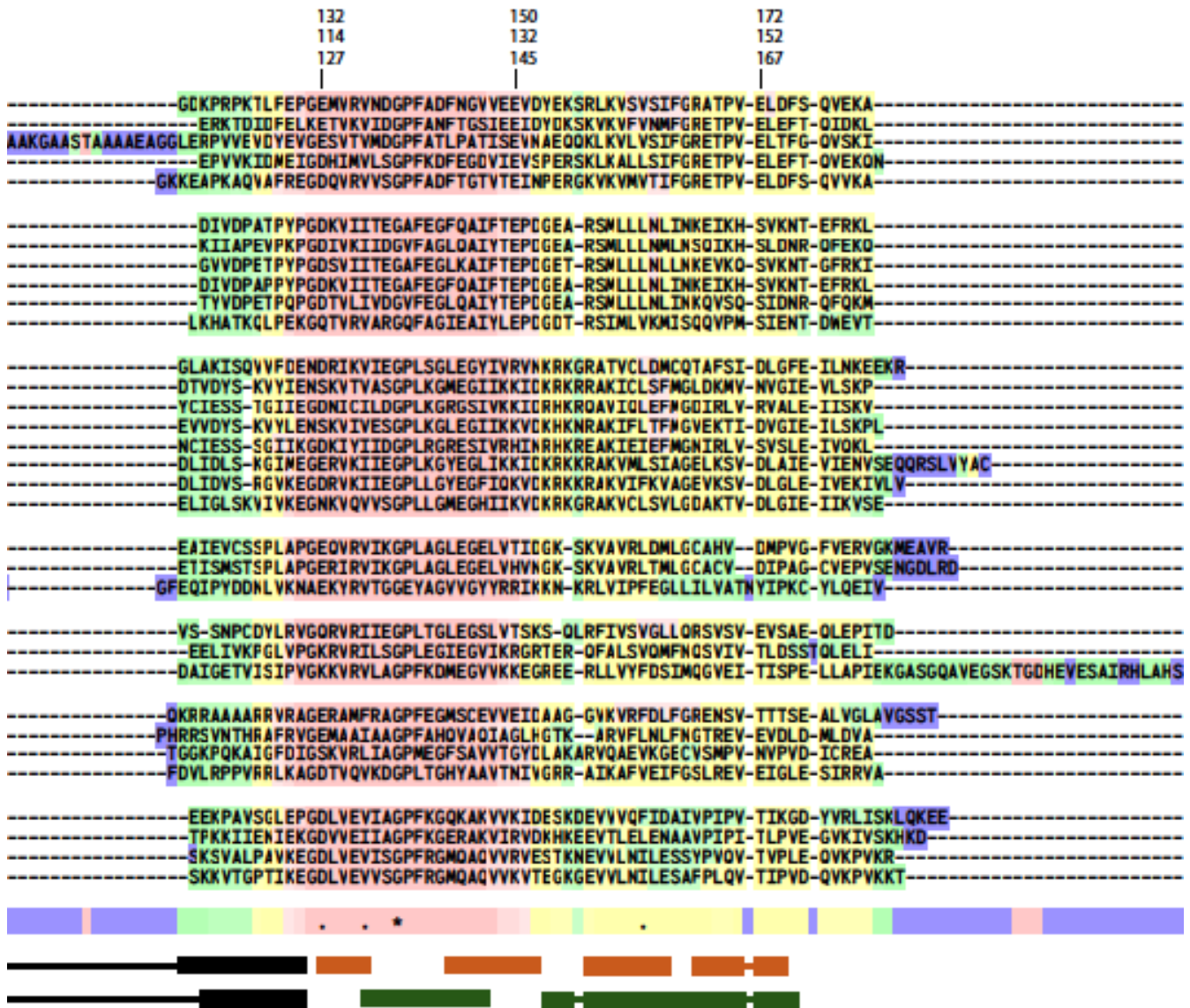
EFFLPTYTPPKS-S--GV--KAKLPLFPGYLFCRYOPLNP--Y--RIVRAPGVIRLLG--GD-----AGPEAVPAQEL EAIRR-----VADSG--
SYLPLKYRRT R--VGGLKRIR--EVEVPLFSGYICFALDREKH--Y--LLYDTQKFVRIIK--VE-----DOSAFVG--ELDAVAK-----AIASG--
SSFLPLIETLRQ--MSDRKK--RVEEPLIRGYVVFVNIYHKEH--Y--HVLETDG VVKFIG--IG-----KTPSVISERDIOWLKR-----LAHEP--

YAFHPVTSRRT R-VR--GKLR--EYERRYLPGYV FARFDGIPIP-HR--VL-TSPFLT GALT--RSD-----GOMGVLGPKRLAALHEMRARDLRQEDERCE--
QVFCPTETKFRKTIK--KRRYSIPVLYPMFCGYIFVGRFSML--ELMAENYITA AVGFDP--EF-----GRRRPAPISDYEMAKLRE-----MSGGLI--
LAYLPMRPGRRR--QPRCKKMI--DNSQLIRGYLFCVCTDFSTGSSVD--EILSCGC VSGLLSFRAD-----KYFHRVPSARVVDIID--HCEQLQ--
TVYCPMLRRETR-HFQSKKML--MKECPLFTGYV FAYLRISD-F--G--TLREMRHVL SVLA--DAG-----GTPIPVAGNIVEDIRD-----AQERGD--

TYNLP-----IYA--ILAPSRVKG YIFVEAPNKG VY-DE--AIRGIRHARGVLP-----GEVPFKEIEHFL-----
KEQLD-----VYS--ILASESLKGYV LVEAETKGOV-EE--LXGMPRVRGIVP-----GTIAIEEIEPLL-----
TNNIN-----EIFS--VVVPPNIKGYVILEATGPHVV-KL--ISSGIRHVKGVAH-----GLIQKEDYTKFV-----
TNNIP-----EIYS--IIVPPGLKGYVIVEASGPHVV-KL--LIAGIRHVRGIAQ-----GLVPKDHIVKMY-----

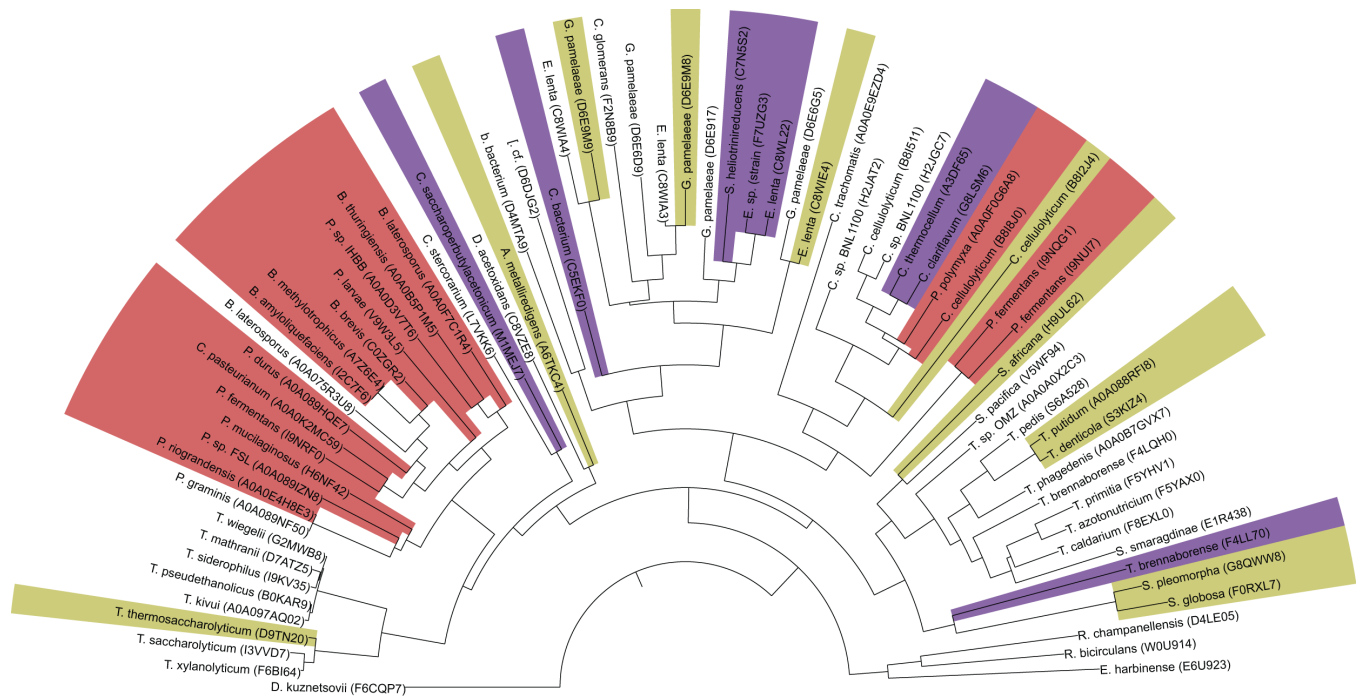
```



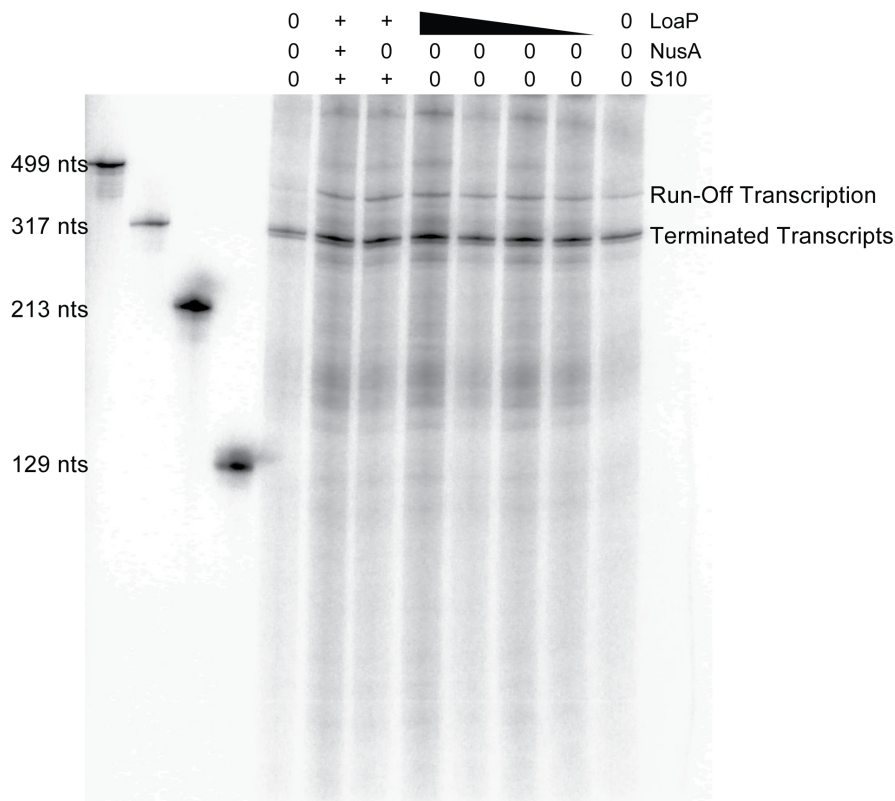


Supplementary Figure 5. Multiple sequence alignment of NusG family protein sequences reveals conserved differences between subtypes of specialized paralogs. Alignment contains 33 protein sequences containing representatives of each subtype of NusG paralog (excluding ActX) aligned by T-COFFEE using the “accurate” alignment method combining PSI-COFFEE and EXPRESSO alignment methods utilizing available crystal structure data. Residues are colored according to T-COFFEE consistency score representing pairwise reliability in alignment with all other sequences. Also shown are secondary structure diagrams representing the secondary structure from PDB structures of *E. coli* NusG

and RfaH (Orange represents beta-strands and green represents alpha-helices, black represents relatively unstructured inter-domain linker).



Supplementary Figure 6. LoaP represents a distinct group of NusG specialized paralogs and is commonly associated with large biosynthetic gene clusters. Subset of large-scale phylogenetic analysis showing LoaP homolog sequences. Background shading represents association of gene sequences with large gene clusters. Red shading denotes LoaP near PKS (polyketide synthase) or NRPS (nonribosomal peptide synthase) gene clusters. Purple shading denotes LoaP near polysaccharide gene clusters. Olive shading denotes LoaP near other types of antiSMASH gene clusters. Unlabeled sequences were not found nearby an antiSMASH predicted gene cluster, although some appear to be next to long stretches of coding sequences in one direction.



Supplementary Figure 7. Transcription of the *dfnA* leader sequence in vitro. The *B. amyloliquefaciens* *dfnA* leader sequence, including the endogenous promoter, was PCR amplified using a reverse oligonucleotide primer that would be predicted to result in either a 314-nucleotide transcript (for premature transcription termination at the putative intrinsic termination site) or a 420-nucleotide transcript (for run-off transcription). Reactions included 2 pmol PCR-generated DNA template, 20 mM Tris-HCl pH 8.0, 15 mM NaCl, 4 mM MgCl₂, 0.1 mM EDTA, 5 mM DTT, 0.01% Triton X100, 1 mM NTPs, 1.5 pmol (5 μ Ci) α -³²P-UTP, and 160 nM σ^A -saturated *B. subtilis* RNAP. LoaP was added to final concentrations of 300 nM, 150 nM, 75 nM, and 3.7 nM. Purified hexahistidine-tagged *B. subtilis* S10 was added to 150 nM and purified hexahistidine-tagged *B. subtilis* NusA was added to 50 nM, as indicated. The reactions were incubated at 37 °C for 1 hour and products were resolved alongside *dfnA*

size markers by 6% urea-denaturing polyacrylamide gel electrophoresis. Gel image is representative of multiple transcription experiments.