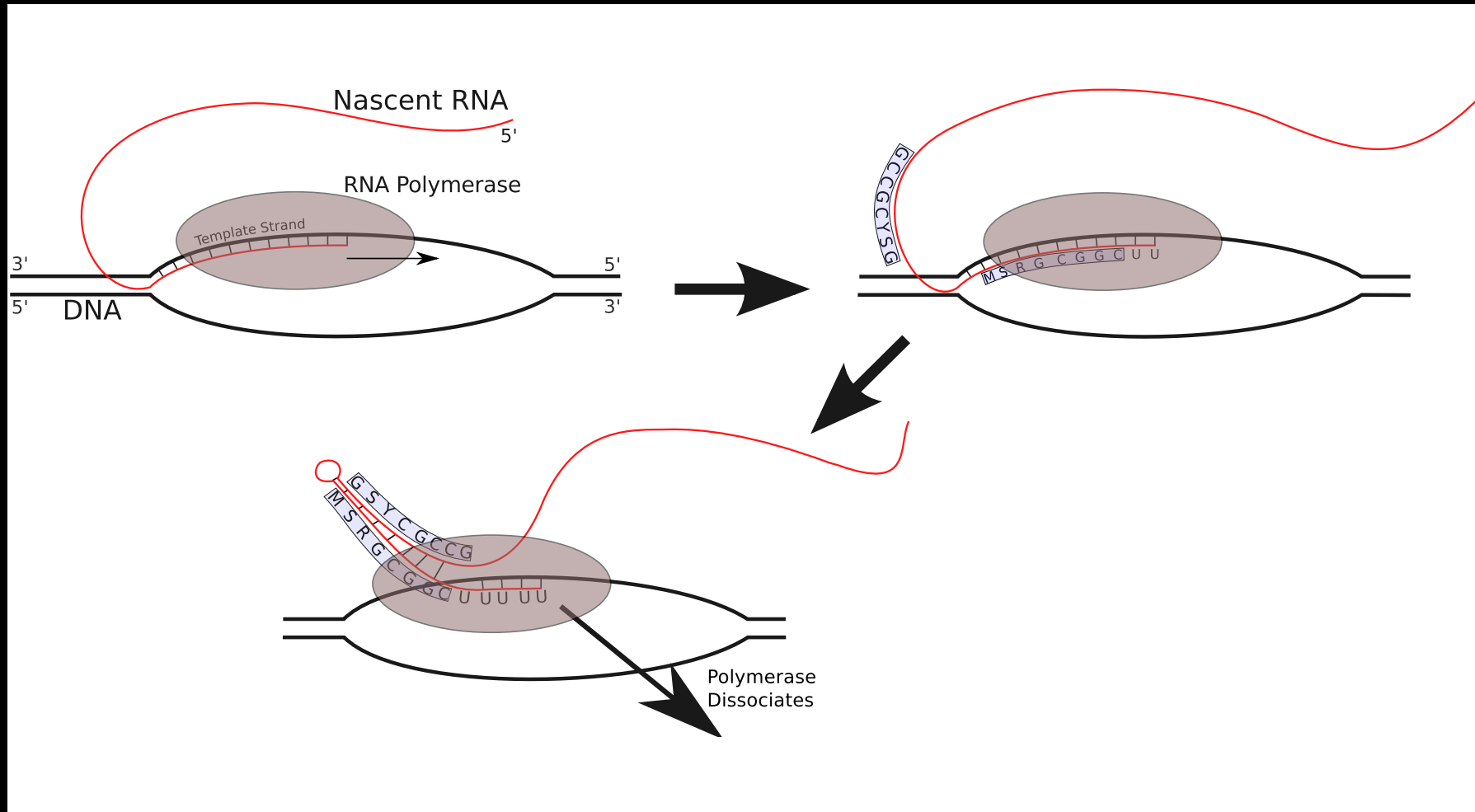


Discovering bacterial
transcription termination
associated motifs through whole
genome clustering

Lars Barquist

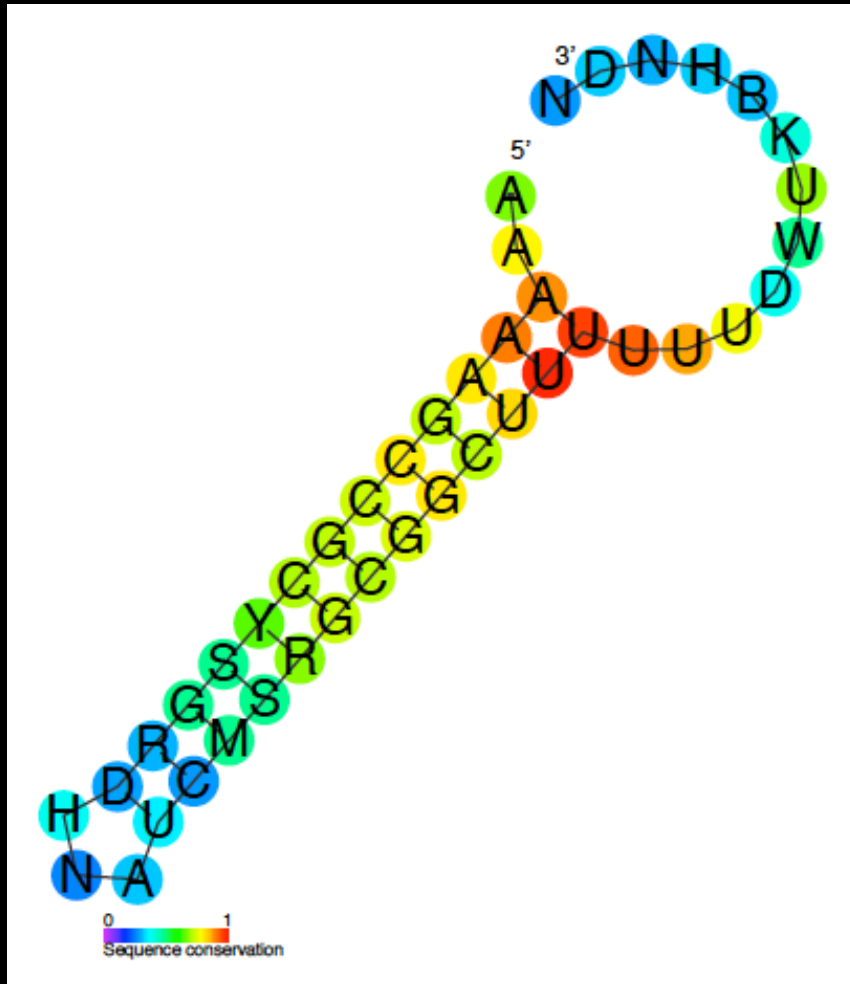
Rho-independent termination



Why care?

- Genomic “landmarks”
- Cis-regulatory components
- Fundamental biological process

A CM for Rho-independent terminators



Constructed from an alignment of 171 *E. coli* and 891 *B. subtilis* experimentally validated terminator sequences.

Published online 7 April 2011

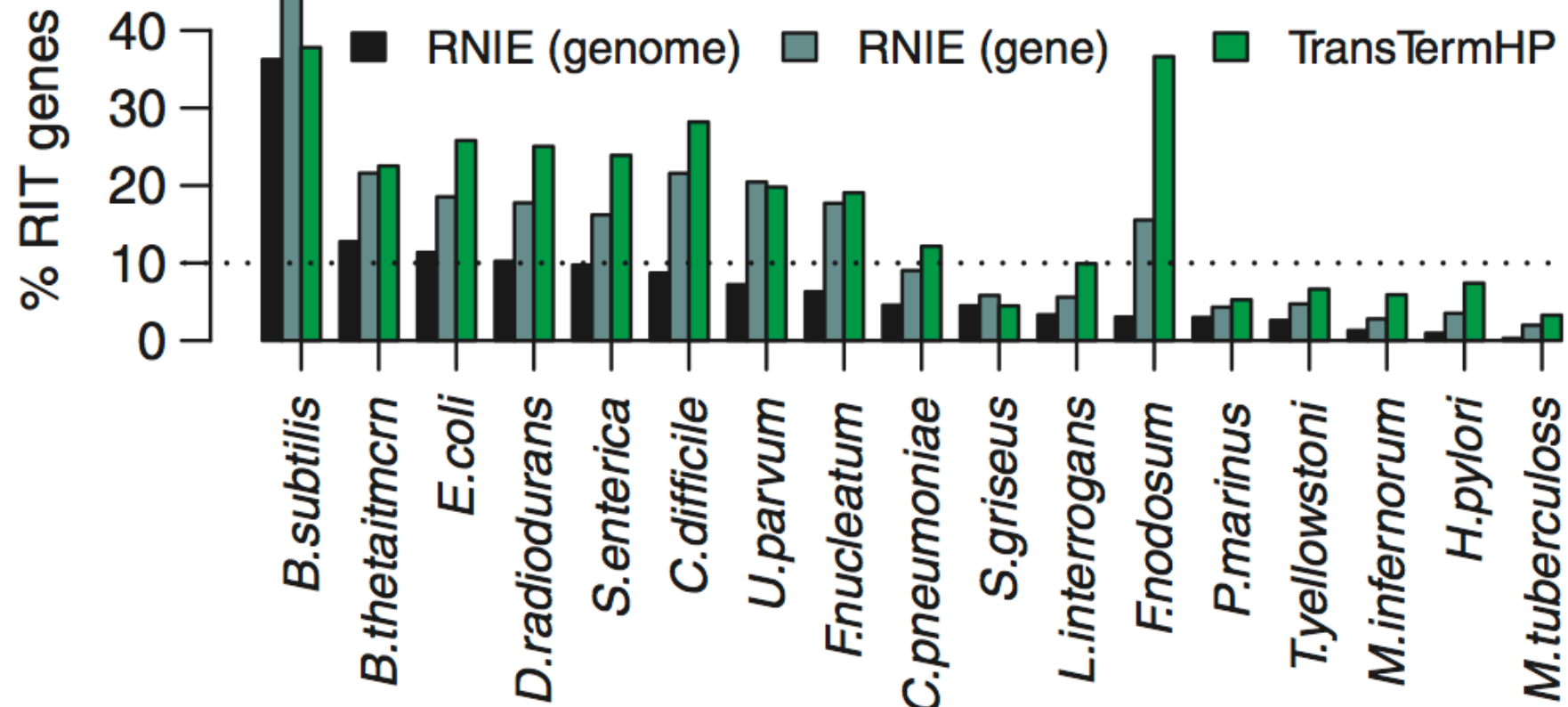
Nucleic Acids Research, 2011, Vol. 39, No. 14 5845–5852
doi:10.1093/nar/gkr168

RNIE: genome-wide prediction of bacterial intrinsic terminators

Paul P. Gardner^{1,*}, Lars Barquist¹, Alex Bateman¹, Eric P. Nawrocki² and Zasha Weinberg³

¹Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Hinxton CB10 1SA0, UK, ²Janelia Farm Research Campus, Howard Hughes Medical Institute, Ashburn, VA 20147 and ³Howard Hughes Medical Institute, Yale University, Box 208103, New Haven, CT 06520, USA

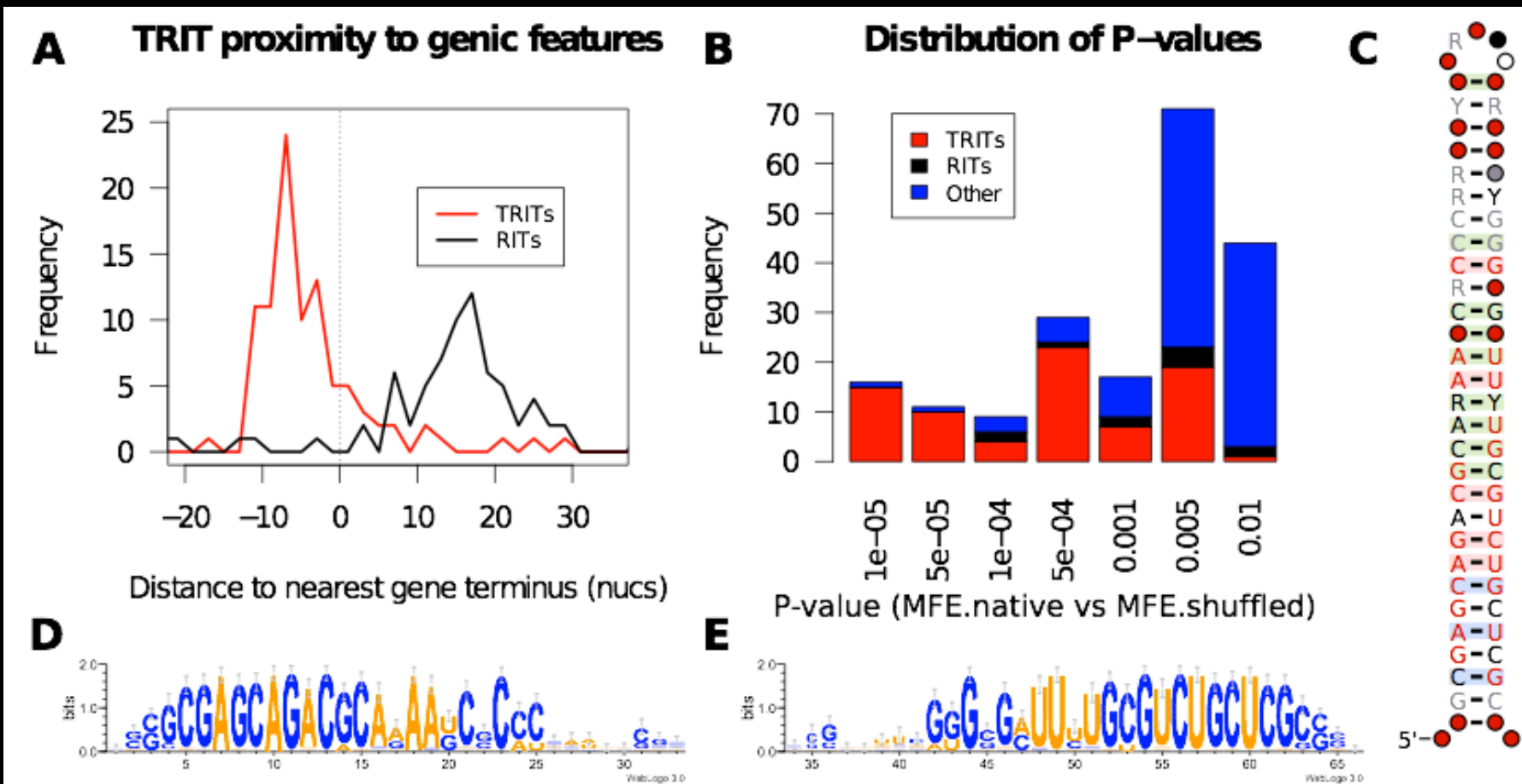
Missing terminators?



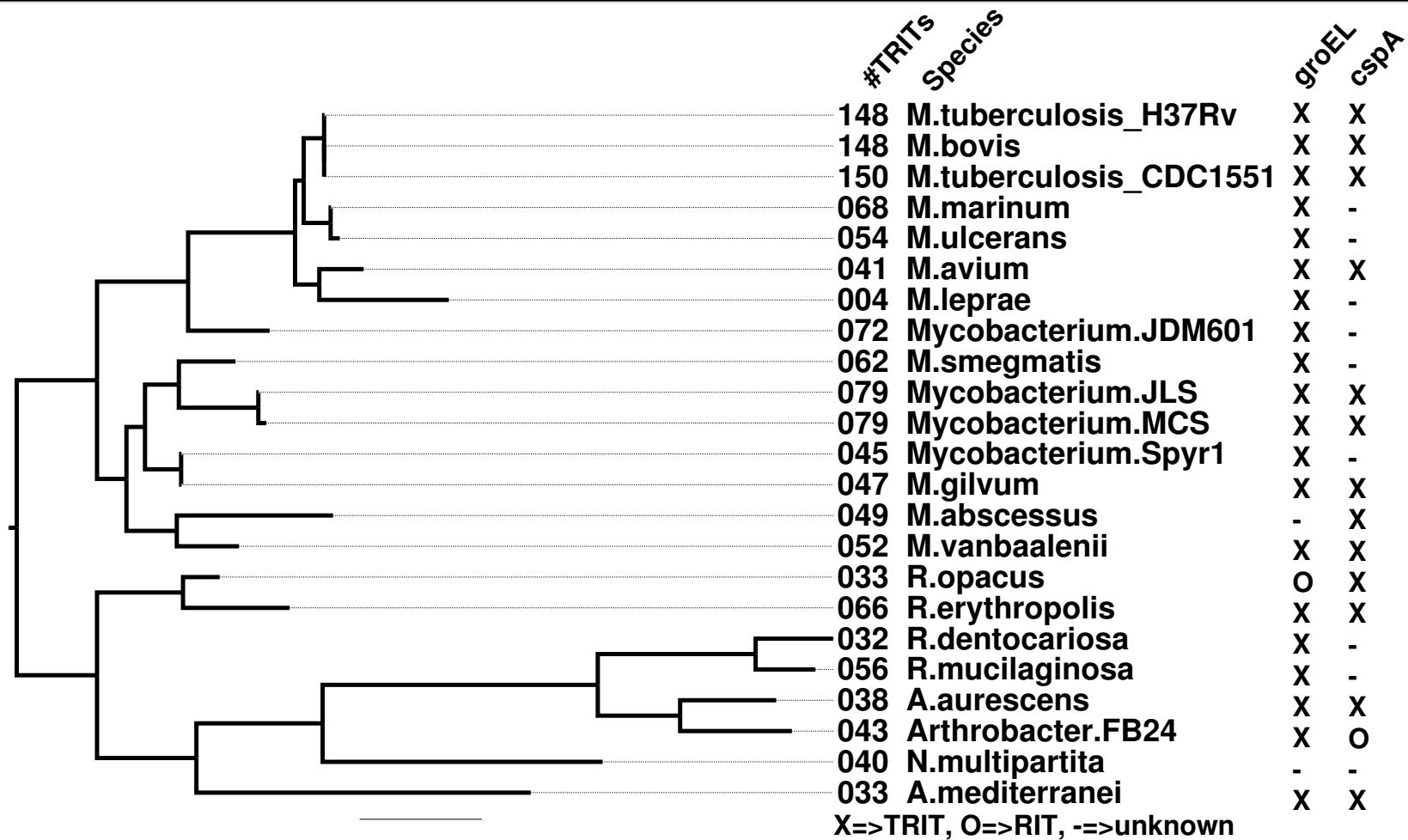
>MT1104|100-bp terminator sequence|1198311,1198410|strand:-1
UGAUCGAGCGCCUCGCGUGAGCGAGGGCGCUCGAGUGCGUUGACUCUGCGUCCACCACGCAAAAUGCGGAGUAGGACACGUGGGUGGGCGCAGAGUCAACGU
>MT2818|100-bp terminator sequence|3054490,3054589|strand:1
GGAUGCUGCUGGUGCUGUAGCCCGGCGAGCAGACGCAAAAUCGCCUCAUUUCGGCAGCAAAAUGGGCGAUUUUGCGUCUGCUCGGCGGGCUACUCGCCGCC
>MT0690|100-bp terminator sequence|757228,757327|strand:-1
GCCUGCUCGGCGACCGCUGACGGCCCCGGGCGAGCAGUCGCAUAAGCCCCCGACACGCCGAGCGUGCGGGGGCUUAUGCGUCUGCUCGCCGGCCGACAAU
>MT0884|100-bp terminator sequence|958341,958440|strand:-1
UGCUGGGCCCCGCAAUUUAGGUUGCGCGAGCAGACGUAAAAGCCCCCGACACGCCGRGCGUGCGGGGGCUUUUACGUCUGCUCGCGCUCAGCUUAGCAGC
>MT2121|100-bp terminator sequence|2319422,2319521|strand:-1
AGGGCGACCUGGAAGCCUGAGUACGGCGAGCAGACGCAGAAUCGCAUGCAGCGGGACCCCGCGCAGUGCGAUUCUGCGUCUGCUCGGCGGGUAGGUUGGCA
>MT1437|100-bp terminator sequence|1567875,1567974|strand:1
ACCUCAAGCGCGCAUCUAGCGUCGAGGGCGCGAGCAGACGCAGAAUCGCACGCGGAAAGGCUUCCGCGUGCGAUUCUGCGUCUGCUCGGCGCUAGCUCG
>MT1127|100-bp terminator sequence|1223826,1223925|strand:1
AGAUCACCGGGCCGCGCUGAGUGCGCCUCCCGGAGCAGACACAGAAUCGCACUGCGCCGGCCGGCGCGUGCGAUUCUGUGUCUGCUCGCCGGUAGACU
>MT0244|100-bp terminator sequence|279624,279723|strand:1
UGGCAGCGGCCUCGGCCUAGGCCUGGCGAGCAGACGCAAAAUCGCCCAAUUUCGUGCCGAAUUGGGCGAUUUUGCGUCUGCUCGCCAGGGGAACGCUAGG
>MT2819|100-bp terminator sequence|3054498,3054597|strand:-1
GGGGUGACGGCGGCGAGUAGCCCGCCGAGCAGACGCAAAAUCGCCCAAUUUCGUGCCGAAUUGAGGCGAUUUUGCGUCUGCUCGCCGGGGCUACAGCACCAG
>MT3684|100-bp terminator sequence|4013504,4013603|strand:1
CCCAGGUUCUGGGGAUCUAGCCGCAAGGGCGCGAGCAGACGCAGAAUCGCAUGAUUUAGAGUCUAAAUCAUGCGAUUCUGCGUCUGCUCGGCAGGCUCGCG
>MT1438|100-bp terminator sequence|1567887,1567986|strand:-1
ACUACCGAUCAGCAGCUAGCGCCGAGCAGACGCAGAAUCGCACGCGGAAGCCUUUCCGCGUGCGAUUCUGCGUCUGCUCGGCGCCUCGACGCUAGAUG
>MT0208|100-bp terminator sequence|234550,234649|strand:-1
GGGUCCGGAUCUGGAACUAGGCUAAGACCCCGGAGCAGACGCAGAAUCGCCCAAUUUGGGCUACCAAUUGGGCGAUUCUGCGUCUGCUCGGCGGAUGGCG
>MT2865|100-bp terminator sequence|3099695,3099794|strand:-1
AACCUCGCCGACCCCGUAGCAGUUGGGCUAGGUUGGCCGAGCAGACGCAAAAAGCCCCCAUUUCGGGCCGAAAUGGGGGCUUUUGCGUCUGCUCAGGCC
>MT3420|100-bp terminator sequence|3702432,3702531|strand:1
CGCUGAUGUUCACCCGUGAGGGCUUGCGCGAGCAGACGCAAAAUCGCCGAAAACCAGUGGUUUUGGGCGAUUUUGCGUCUGCUCGGCGAGCCGGGUCU
>MT0670|100-bp terminator sequence|737658,737757|strand:-1
AGCCGGGUGCCGCGGCCUAAGCCCGGCGAGCAGACGCAUAAGCCCCCGACACGCCGUGCGUGCGGGGGCUUAUGCGUCGCCGGGGAAACUACUCCCCCG
>MT0669.1|100-bp terminator sequence|737650,737749|strand:1
GCAACUUCGCGGGGAGUAGUUUCCCCGCGACGCAUAAGCCCCCGCACGCACGGCGUGUCGGGGCUUAUGCGUCUGCUCGCCGGGCUUAGGCCGCGGC
>MT3441|100-bp terminator sequence|3722137,3722236|strand:1
CGCAUGUCCGGGCCAGUUAGCGGGCGGAGCAGACGCAAAAUCGCCCAAUUUCGGCAGCAAAAUGGGCGAUUUUGCGUCUGCUCGCCCUAAUUGGCCAGCU
>MT0245|100-bp terminator sequence|279639,279738|strand:-1
AGACCGUCUGGAUCGCCUAGCGUUCUCCUGGCGAGCAGACGCAAAAUCGCCCAAUUUCGGCACGAAAUUGGGCGAUUUUGCGUCUGCUCGCCAGGCCUAGG
>MT0883|100-bp terminator sequence|958334,958433|strand:1
CGGAGUCGUCUAAGCUGAGCGCGAGCAGACGUAAAAGCCCCCGCACGCYCGGCGUGUCGGGGGCUUUUACGUCUGCUCGCGCAACC UAAAUUGCCGGG

>MT1104|100-bp terminator sequence|1198311,1198410|strand:-1
UGAUCGAGCGCCUCGCGUGAGCGAGGGCGCUCGAGUGCGUUGACUCUGCGUCCACCACGCAAAAUGCGAGUAGGACACGUGGGUGGGCGCAGAGUCAACGU
>MT2818|100-bp terminator sequence|3054490,3054589|strand:1
GGAUGCUGCUGGUGCUGUAGCCCG**GCGAGCAGACGCAAAU**CGCCUCAUUUCGGCACGAAAUG**GGGCGAUUUUGCGUCUCUGCG**CGGGCUACUCGCCGCC
>MT0690|100-bp terminator sequence|757228,757327|strand:-1
GCCUGCUCGGCGACCGCUGACGGCCCCGG**GCGAGCAGUCGCAUAAGCCCCCG**ACACGCCGAGCGUG**CGGGGGCUUAUGCGUCUCUCGC**CGGCCGACAAU
>MT0884|100-bp terminator sequence|958341,958440|strand:-1
UGCUGGGCCCCGCAAUUUAGGUUGC**GCGAGCAGACGUA**AAAAGCCCCGACACGCCGRGCGUG**CGGGGGCUUUUACGUCUCUCGC**GCUCAGCUUAGCAGC
>MT2121|100-bp terminator sequence|2319422,2319521|strand:-1
AGGGCGACCUGGAAGCCUGAGUACG**GCGAGCAGACGCAGAAUCGCAUG**CGCGGGACCCCGCG**CAGUGCGAUUCUGCGUCUCUCGC**CGCGGUAGGUUGGCA
>MT1437|100-bp terminator sequence|1567875,1567974|strand:1
ACCUCAAGCGCGCAUCUAGCGUCGAGG**GCGCGAGCAGACGCAGAAUCGCAC**GCGGAAAGGCUUCCGC**GUGCGAUUCUGCGUCUCUCGGCGC**UAGCUCG
>MT1127|100-bp terminator sequence|1223826,1223925|strand:1
AGAUCACCGGGCCGCGCUGAGUGCGCCUCCCGGAGCAGACACAGAAUCGCACUGCGCCGGCCGGCGGUGCGAUUCUGUGUCUCUCGCGCCGGUAGACU
>MT0244|100-bp terminator sequence|279624,279723|strand:1
UGGCAGCGGCCUCGGCCUAGGCCUG**GCGAGCAGACGCAAAU**CGCCAAUUUCGUGCCGAAUUG**GGGCGAUUUUGCGUCUCUCGC**CAGGGGAACGCUAGG
>MT2819|100-bp terminator sequence|3054498,3054597|strand:-1
GGGGUGACGGCGGCGAGUAGCCC**GCCGAGCAGACGCAAAU**CGCCAAUUUCGUGCCGAAUUG**GGCGAUUUUGCGUCUCUCGC**CGGGCUACAGCACCAG
>MT3684|100-bp terminator sequence|4013504,4013603|strand:1
CCCAGGUUCUGGGGAUCUAGCCGCAAGG**GCGCGAGCAGACGCAGAAUCGCA**UGAUUUUGAGCUCAAAUCA**UGCGAUUCUGCGUCUCUCGC**GAGGCUCGCG
>MT1438|100-bp terminator sequence|1567887,1567986|strand:-1
ACUACCGAUCAGCAGCUAG**GCGCCGAGCAGACGCAGAAUCG**CACGCGGAAGCCUUUCCGCG**UGCGAUUCUGCGUCUCUCGCGC**CCUCGACGCUAGAUG
>MT0208|100-bp terminator sequence|234550,234649|strand:-1
GGGUCCGGAUCUGGAACUAGGCUAAGACCCCGGAGCAGACGCAGAAUCGCCAAUUUGGGCUACCAAUUGGGCGAUUCUGCGUCUCUCGCGCGGAUGGCG
>MT2865|100-bp terminator sequence|3099695,3099794|strand:-1
AACCUCGCCGACCCCGUAGCAGUUGGGCUAGGUUGGCCGAGCAGACGCAAAAGCCCCCAUUUCGGGCCGAAAUGGGGGCUUUUGCGUCUCUCACGCC
>MT3420|100-bp terminator sequence|3702432,3702531|strand:1
CGCUGAUGUUCACCCGUGAGGGCUU**GCGCGAGCAGACGCAAAU**CGCCGAAAACCAGUGGUUUU**GGGCGAUUUUGCGUCUCUCGC**GCAGCCGGGUCU
>MT0670|100-bp terminator sequence|737658,737757|strand:-1
AGCCGGGUGCCGCGGCCUAAGCCCG**GCGAGCAGACGCAUAAGCCCCCG**ACACGCCGUGCGUG**CGGGGGCUUAUGCGUCGC**CGGGGAAACUACUCCCCCG
>MT0669.1|100-bp terminator sequence|737650,737749|strand:1
GCAACUUCGCGGGGAGUAGUUUCCCG**GCGACGCAUAAGCCCCCG**CACGCACGGCGUGU**CGGGGGCUUAUGCGUCUCUCG**CCGGGCCUAGGCCGCGGC
>MT3441|100-bp terminator sequence|3722137,3722236|strand:1
CGCAUGUCCGGGCCAGUUAGCGG**GCGCGAGCAGACGCAAAU**CGCCAAUUUCGGCACGAAAU**UGGGCGAUUUUGCGUCUCUCGC**CCUAAUUGGCCAGCU
>MT0245|100-bp terminator sequence|279639,279738|strand:-1
AGACCGUCUGGAUCGCCUAGCGUUCUCCUG**GCGAGCAGACGCAAAU**CGCCAAUUCGGCACGA**AAUUGGGCGAUUUUGCGUCUCUCGC**CAGGCCUAGG
>MT0883|100-bp terminator sequence|958334,958433|strand:1
CGGAGUCGUCUAAGCUGA**GCGCGAGCAGACGUA**AAAAGCCCCGACGCYCGGCGUGU**CGGGGGCUUUUACGUCUCUCGC**GCAACCUAAAUUGCCGGG

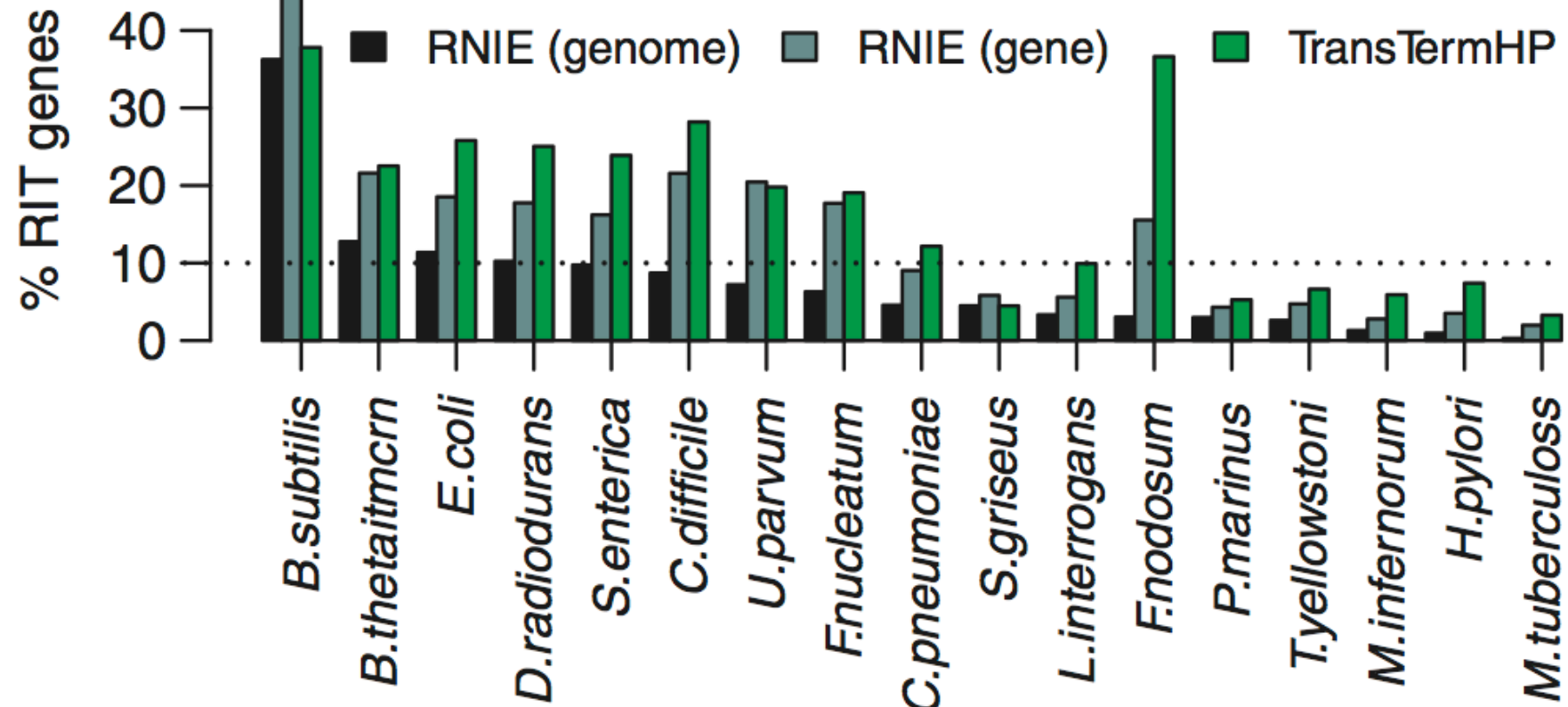
Mycobacterium tuberculosis



Missing terminators?



Missing terminators?



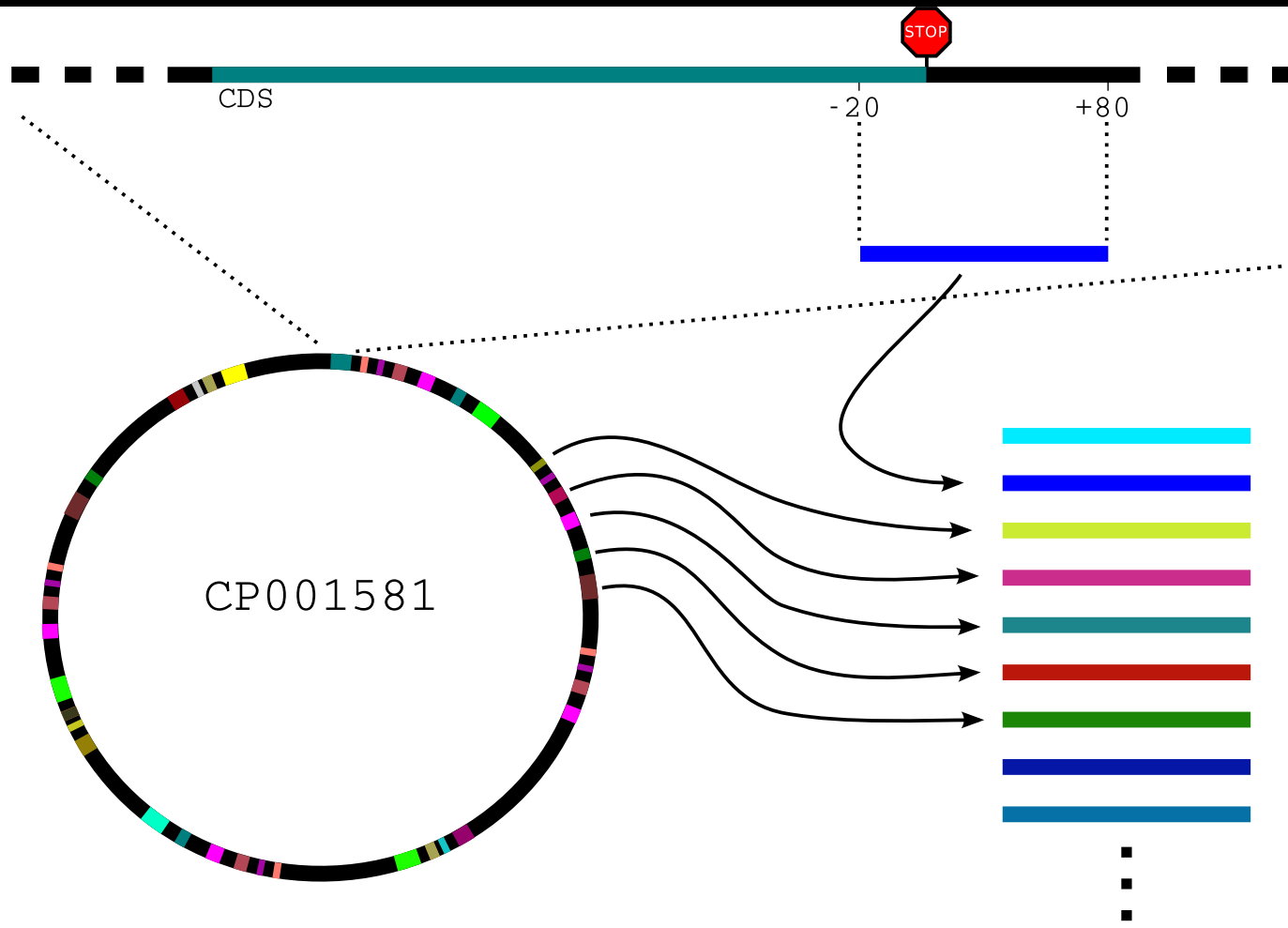
Questions

Are there terminator "families"?

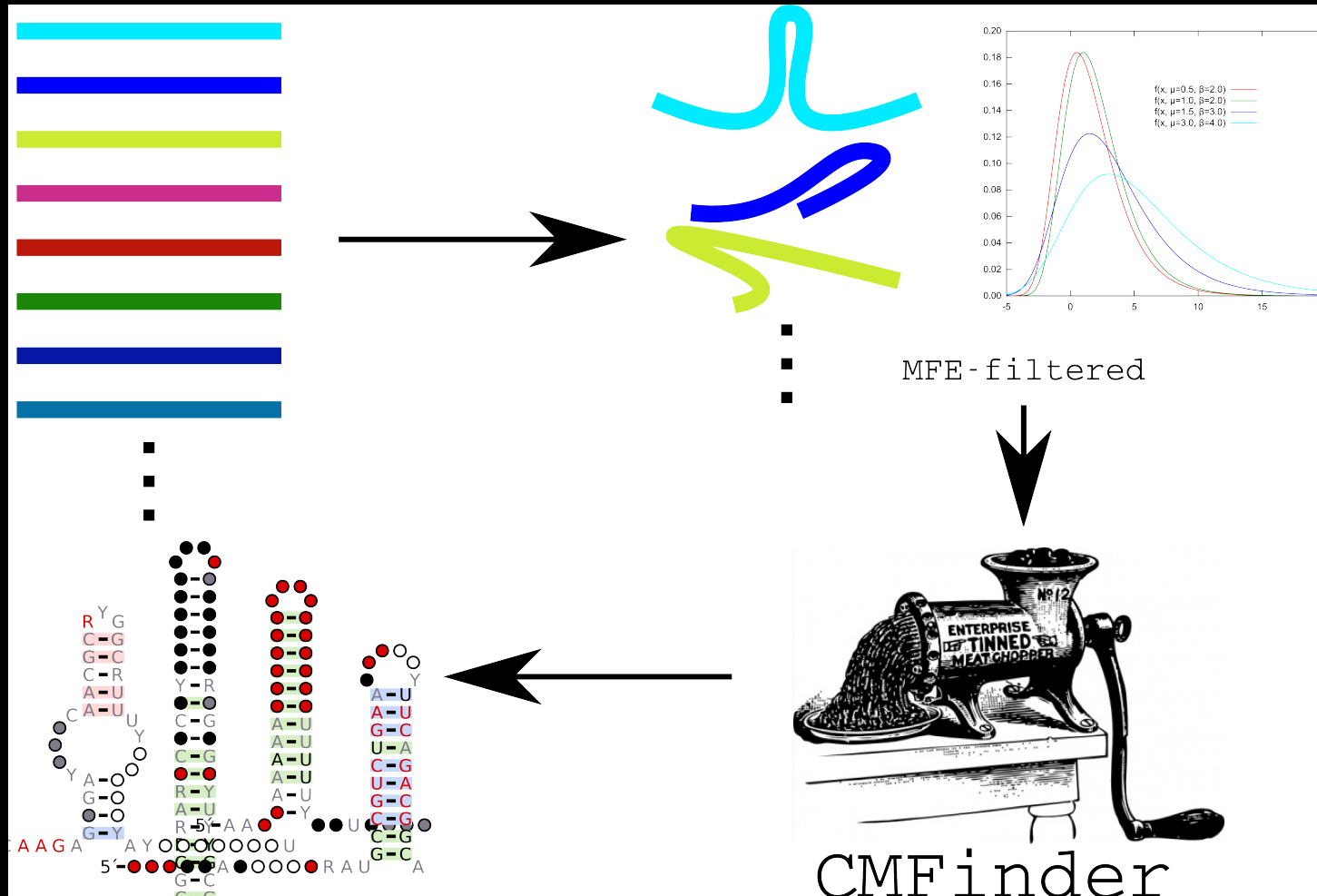
How many are there?

Do these families have a shared
origin?

Mining EMBL

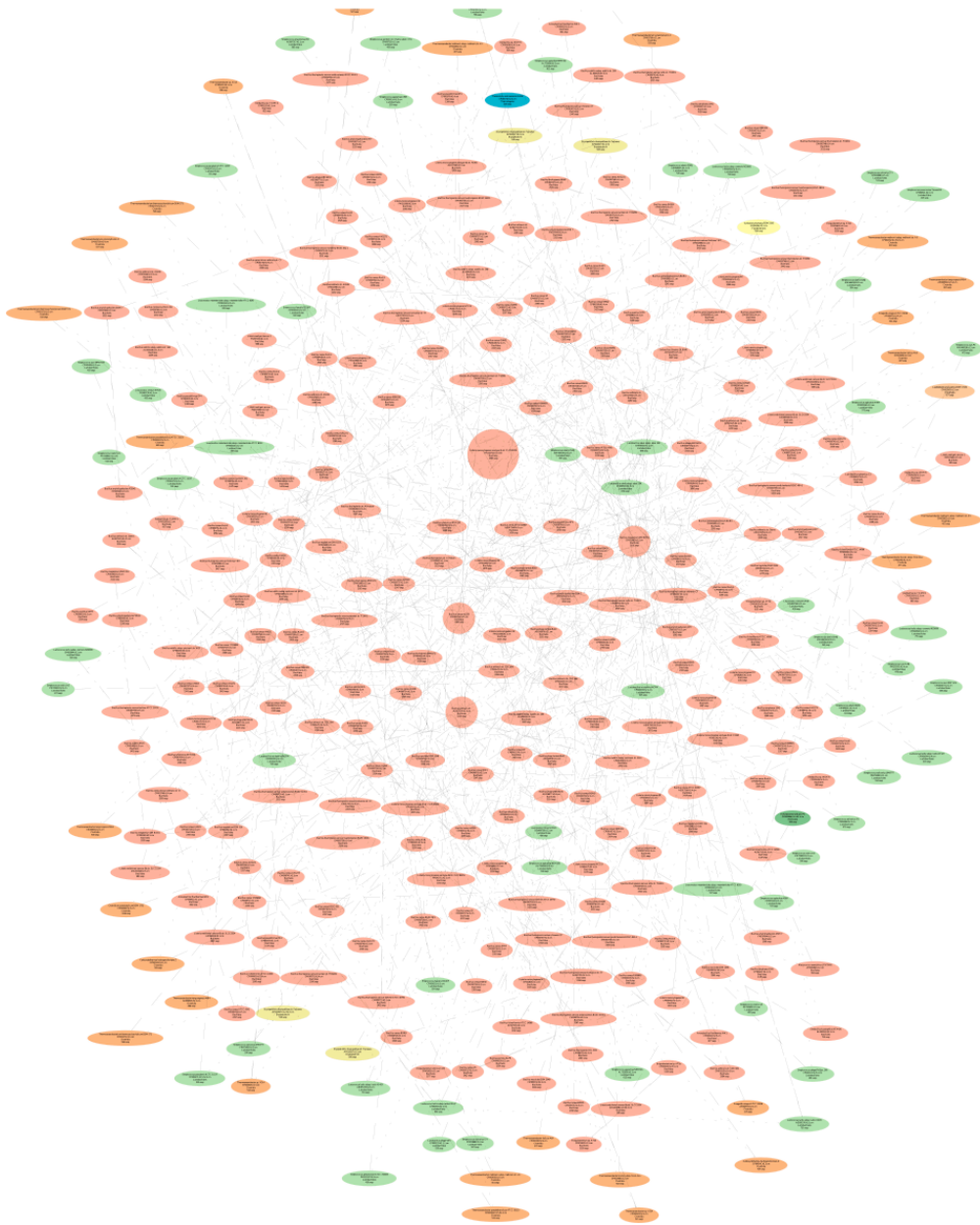


Mining EMBL

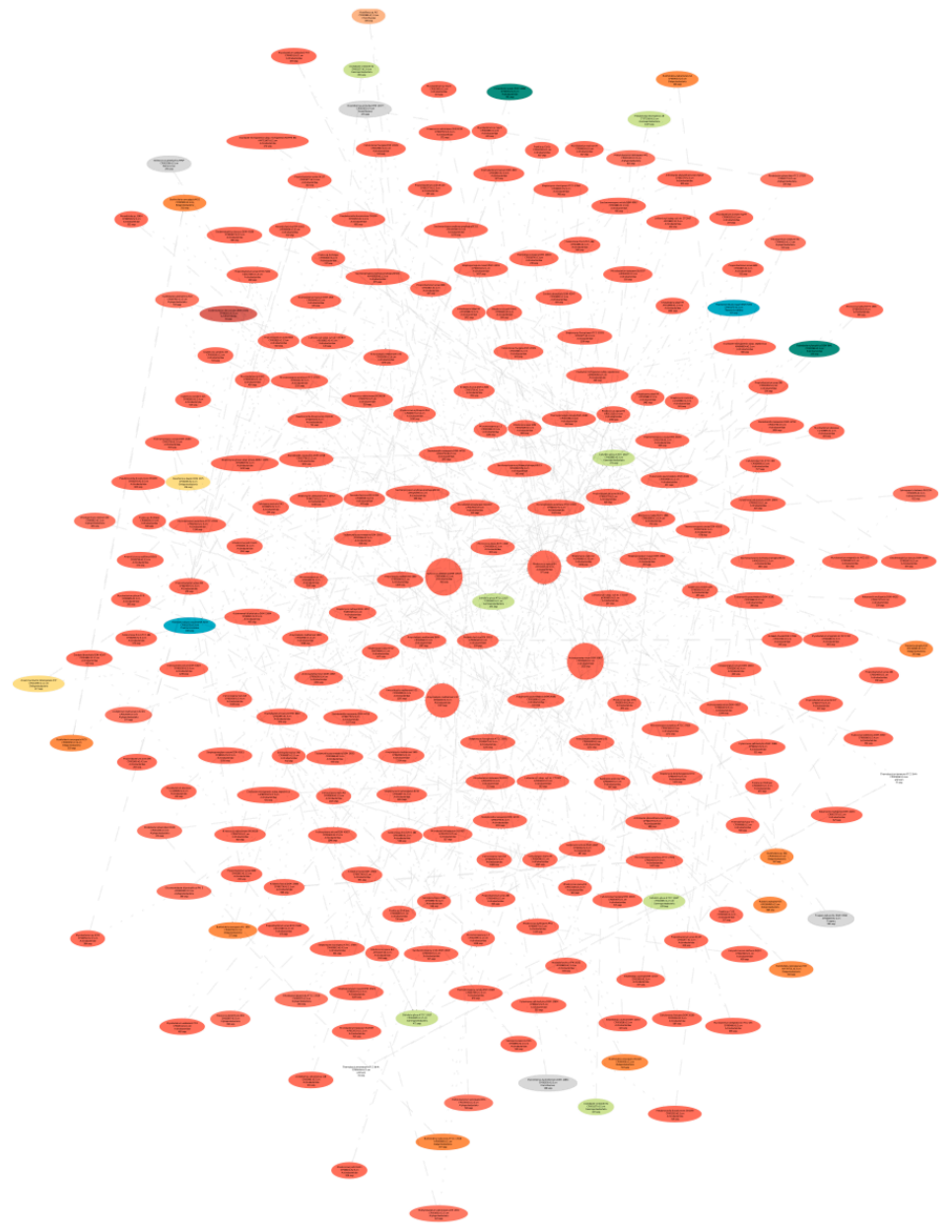


Clustering CMs

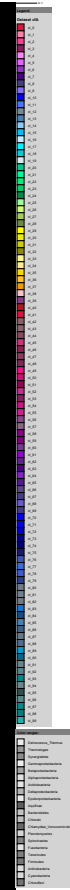
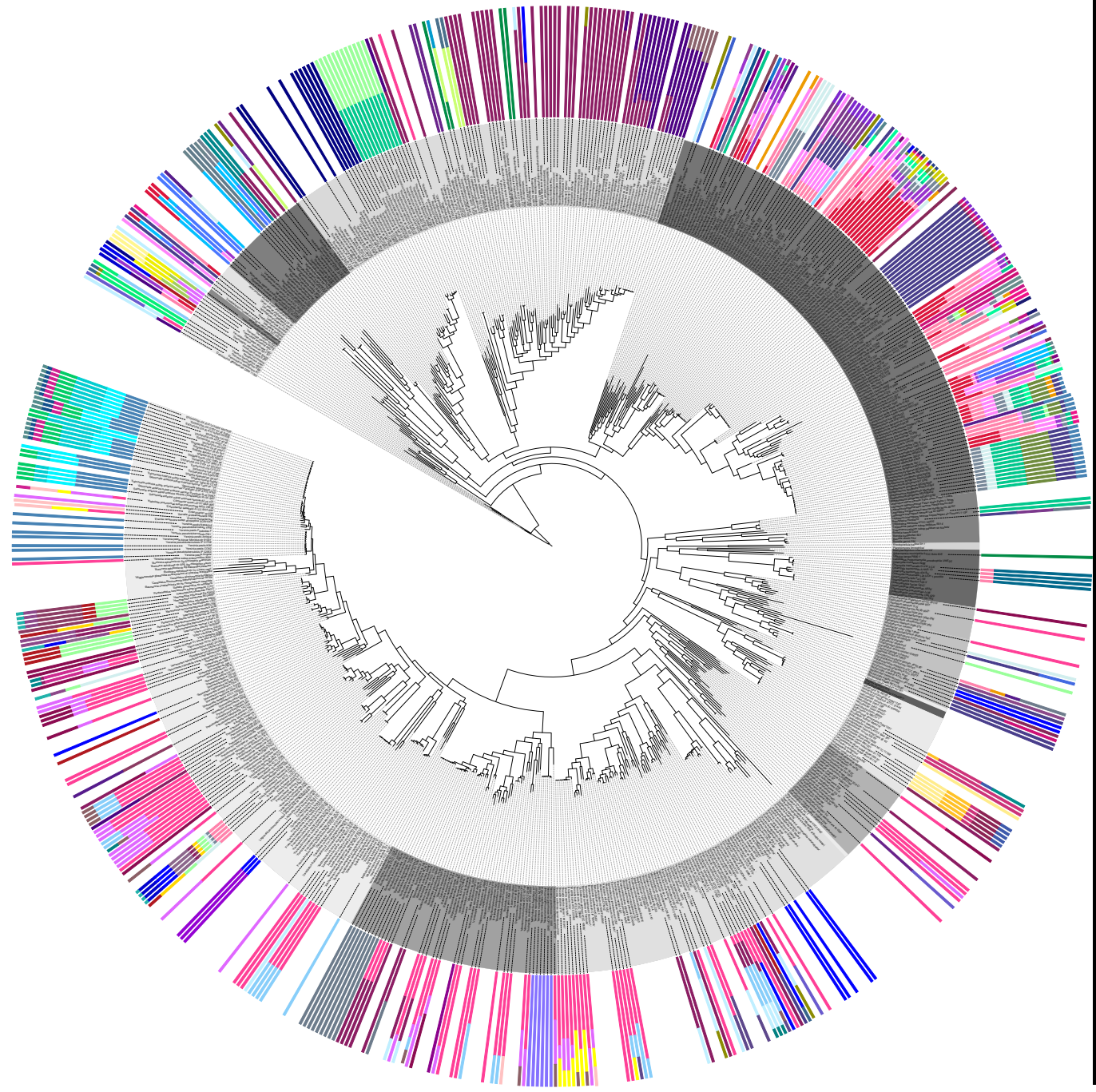
- Emit 1000 sequences from one model, score with the other & vice versa
- Average $-\lg(\text{E-values})$ & shift positive
- Feed in to MCL (see Enright et al. TRIBE-MCL)
- Results in ~100 (non-trivial) clusters



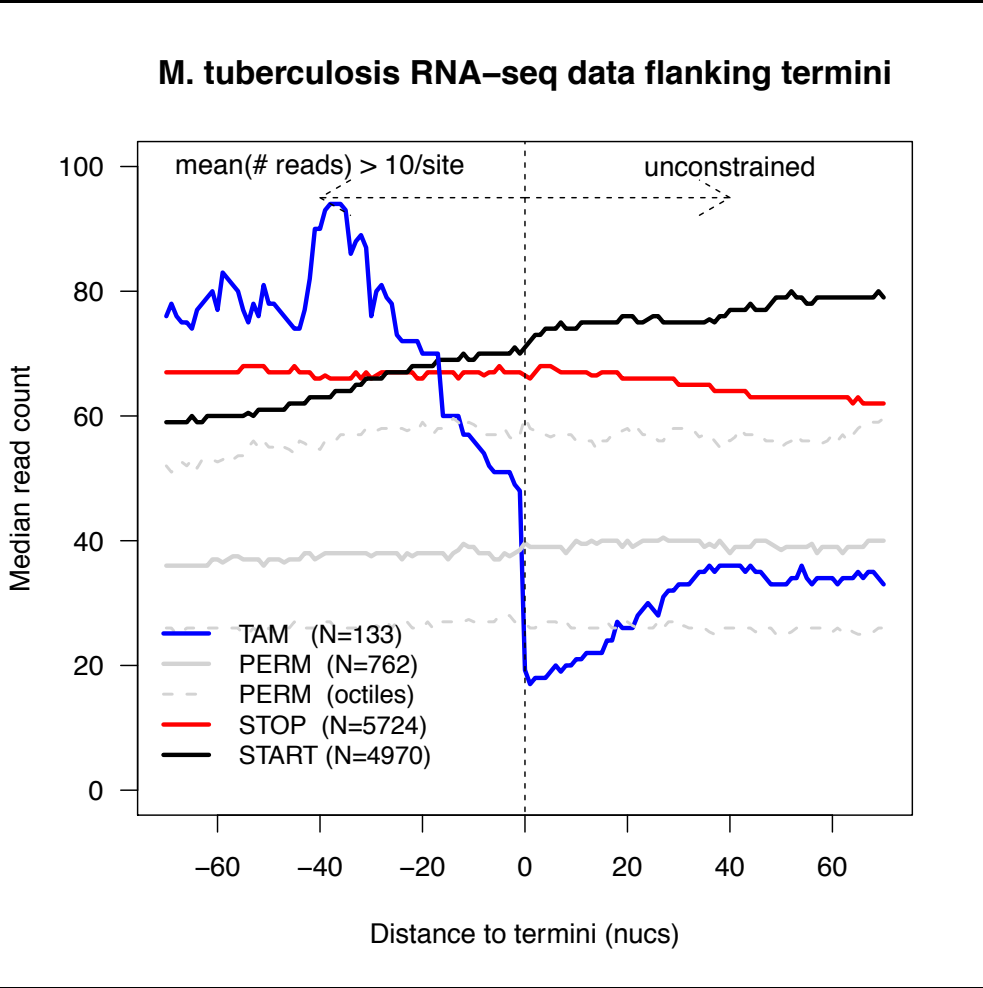
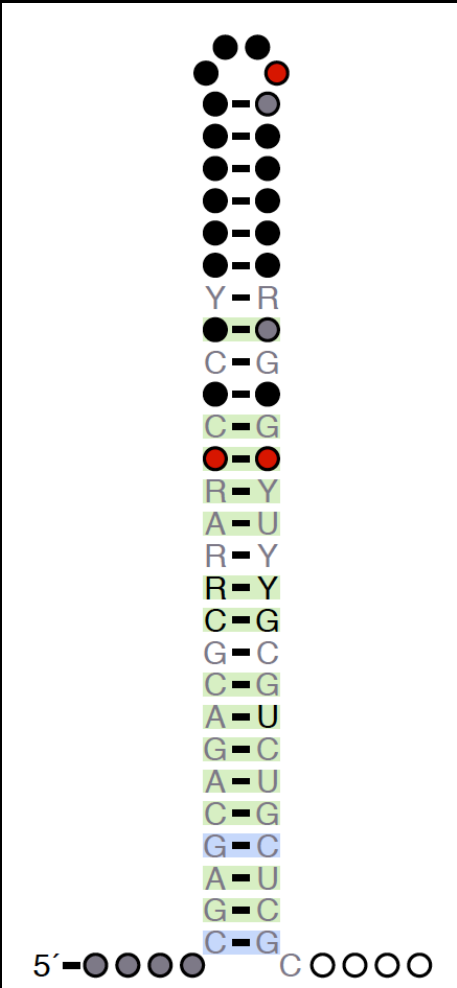
Cluster 0
73% Bacillales



Cluster 3
81% Actinobacteria

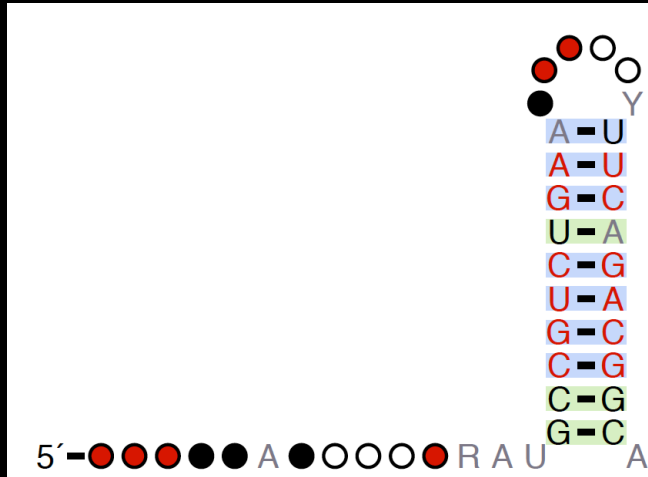


Actinobacterial Terminator



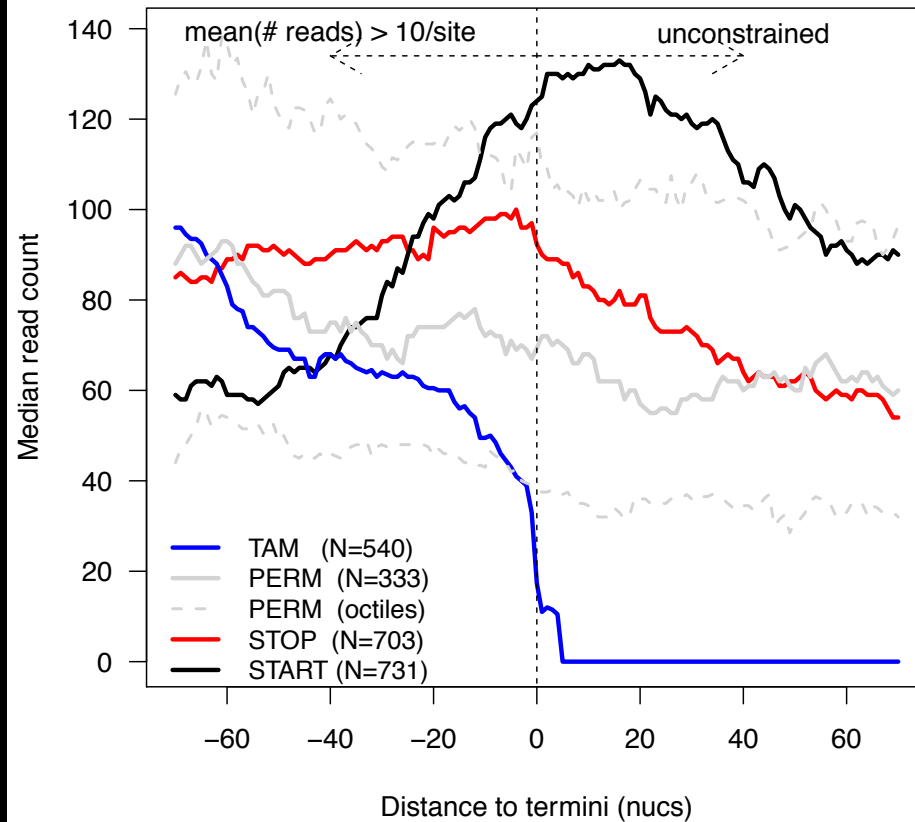
Data courtesy of KB Arnvig, NIMH

Neisseria DNA Uptake Terminator

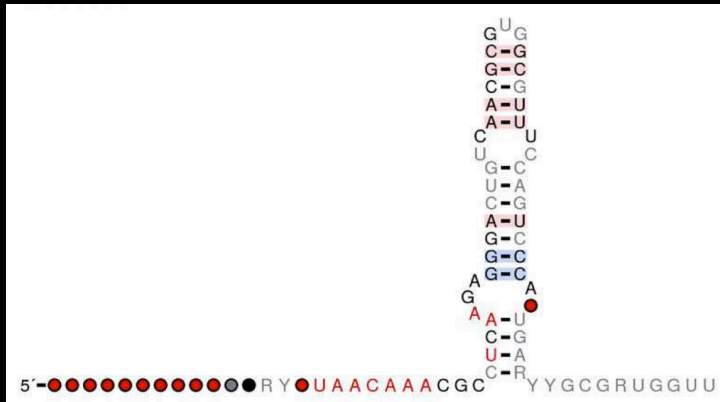


~600/genome, high percent id
 Patchy poly-U tail
 Exapted DNA Uptake sequence?

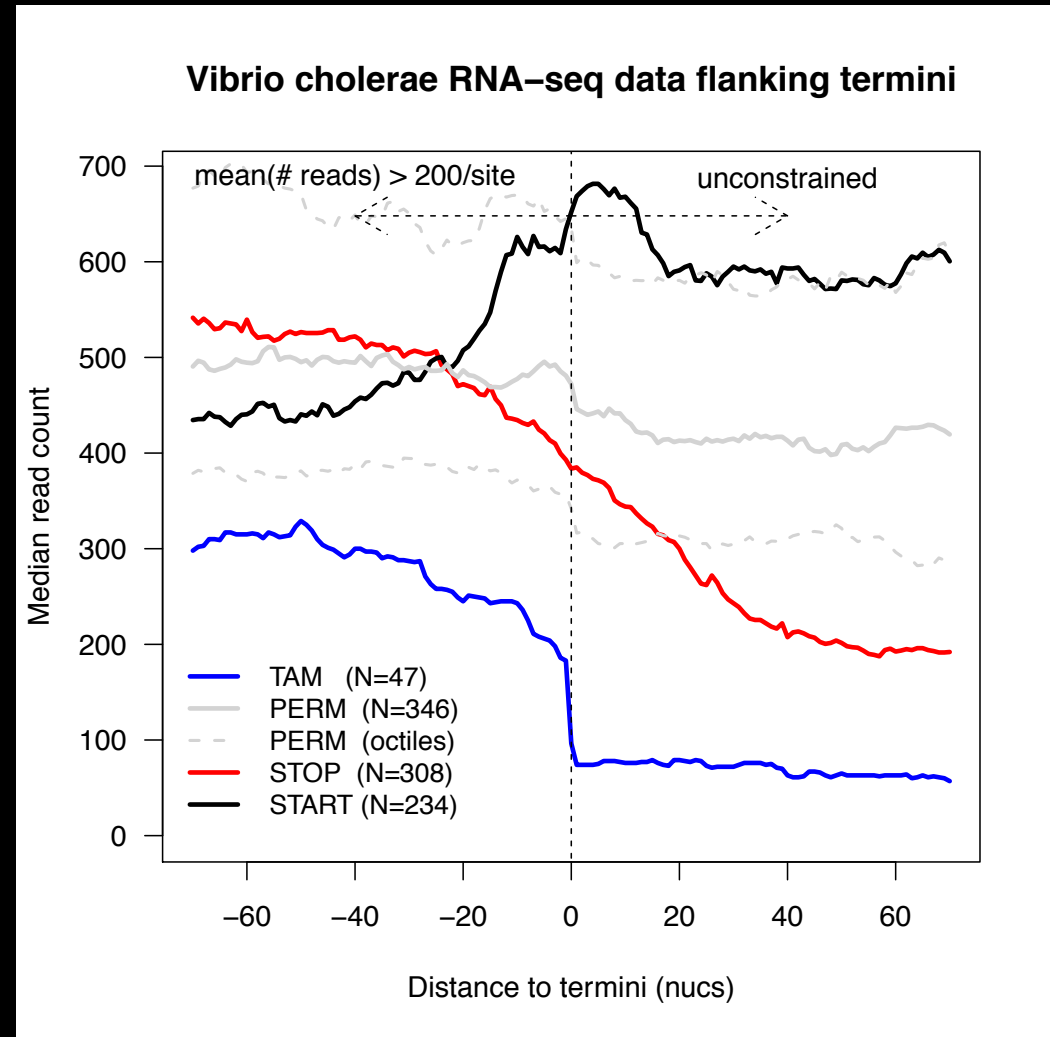
N. gonorrhoeae RNA-seq data flanking termini



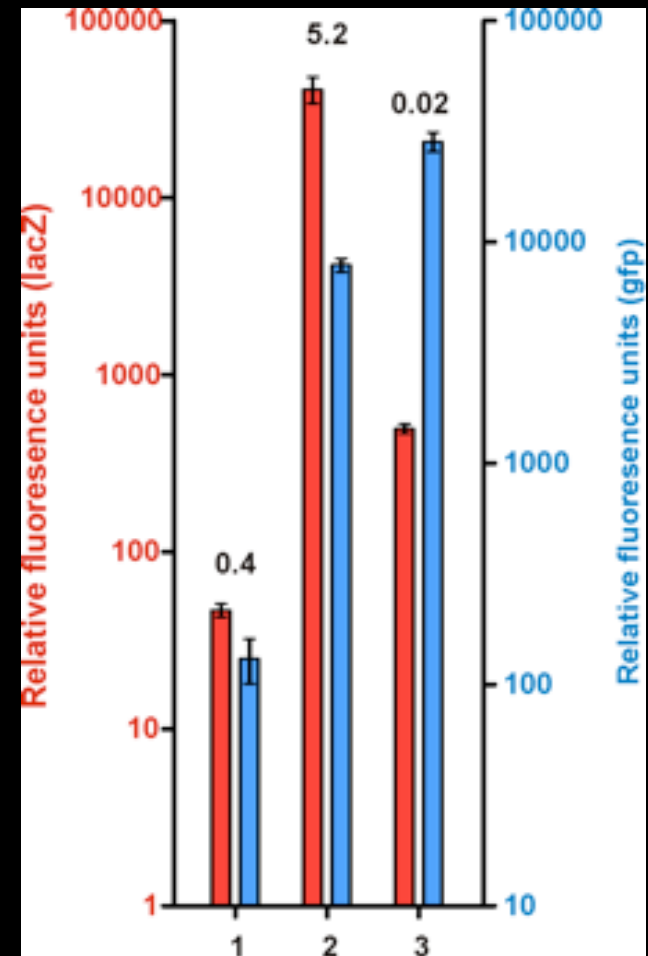
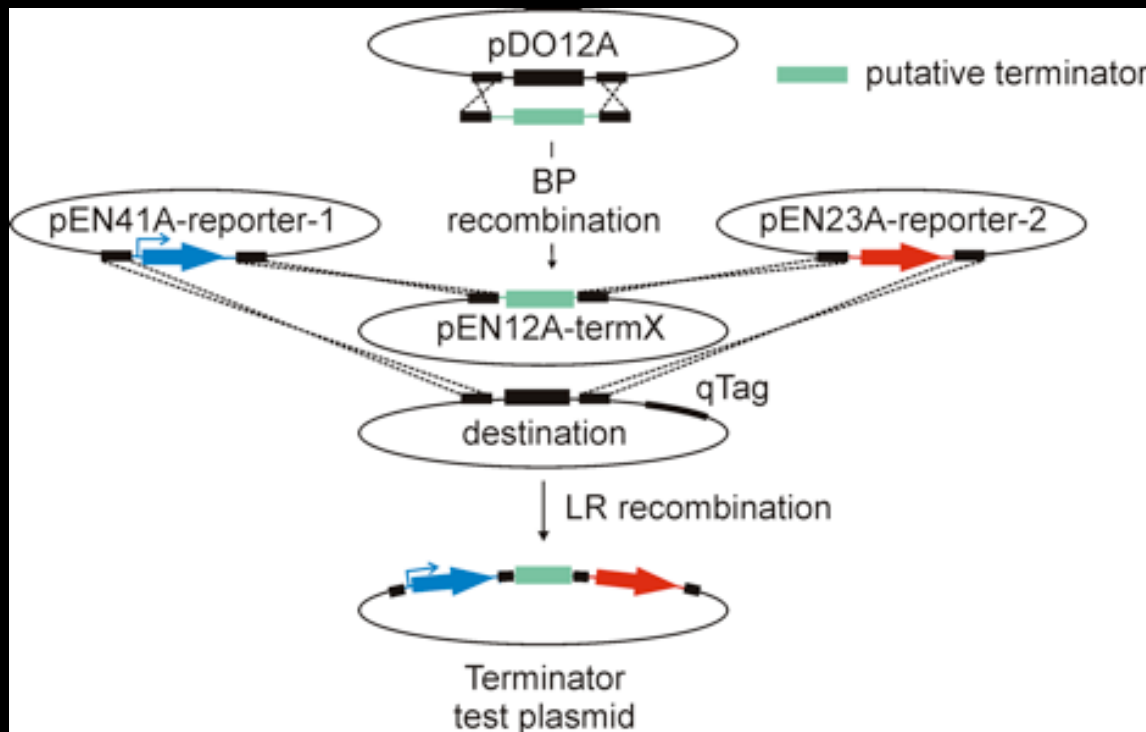
Vibrio Superintegron Terminus-Associated Motif



attC SI recombination site



Experimental Validation



Courtesy of D Schnappinger,
Weill Cornell Medical College

Conclusions so far...

- Rho-independent termination has arisen independently multiple times
- G/C-rich hairpin + poly-U tail appears to be the dominant (but not the only!) model for RI termination
- Connections to horizontal transfer?

Acknowledgements

Paul Gardner University of Canterbury
James Hadfield New Zealand

Tine Arnvig MRC National Institute
for Medical Research, UK

Dirk Schnappinger Weill Cornell
Medical College, USA