Determination of Tandem Repeats in Bacillus Bacteriophages
Ellen Korcovelos

Introduction

Tandem repeats, or highly repetitive DNA sequences, are found in many aspects of biological systems, such as the 3-nt repeat of CAG in Huntington's disease and the poly-A tail of transcriptional termination. These tandem repeats are described as a continuous pattern of nucleic units that are repeated adjacently. Further, while these repeats have been studied greatly in the human genome in association with neurodegenerative diseases[1-3] as well as in bacterial genomes[4-7], there has been little discussion of the purpose or presence of repeats in host bacteria's phage counterparts.

Due to the recombinant nature of prophages, genetic material is readily exchanged between phage and bacteria, and more often than not, novel functions are acquired by the host [8-9]. In some instances, the phages themselves may gain bacterial DNA as a result of illegitimate recombination, in which a novel gene is inserted with no sequence repeats bordering the novel region[10-11]. In cases like that of the Sterne strain in *Bacillus anthracis* bacteria, fused genetic material is capable of altering the function of the cell into that of an innocuous phenotype; in this instance, the toxic components of the anthrax pathogen are protected against by an induced system of antigens. The result of this recombination event is that of immunity against the toxin to the host species (often animals and humans for *B. anthracis*)[12].

Other strains of *Bacillus* bacteria exhibit this same recombinant addition and deletion of protein-coding genes (and thus function), yet study of tandem repeats in phages has only been limited to intramolecular recombination events. In these cases, novel genetic material is inserted into site-specific areas denoted as attachment sites (attP). Serving as markers for these attachment sites, short repeated sequences are often found to flank the inserted prophage DNA[13]. However, a question that does not yet permeate the field of phage genomics is whether tandem repeats exist in phages outside of these attP sites, and if so, if they exist as a result of frequent interactions with their hosts. The present study attempts to identify areas of tandem repeats in twenty-four Bacillus strains. Attention will be especially focused on those repeats found encompassed within genes, which may serve as an indication of phage-host recombination.

Methods

Twenty-four strains of Bacillus phage were run through an autocorrelate function to search for the presence of tandem repeats (Appendix Figure 4). This function, programed through Phantome/BioBIKE[14], consists of a for-each loop that utilizes a shifting frame in order to determine percent correlation in a window of one hundred nucleotides. The offset of the window is set to 7-nt, in which each nucleotide is compared to the seventh nucleotide following it. If the nucleotides match, such as in the sequence **A**TTACAG**A**TTAGAG, the score of the autocorrelate function is incremented by one (Figure 1). The frame is then shifted to the next nucleotide, until all nucleotides are accounted for within the 100-nt window. Once this window is completely checked of all correlated



ATTACAGATTAGAG   score = 1
ATTACAGATTAGAG   score = 2
ATTACAGATTAGAG   score = 3

Window Size = 14

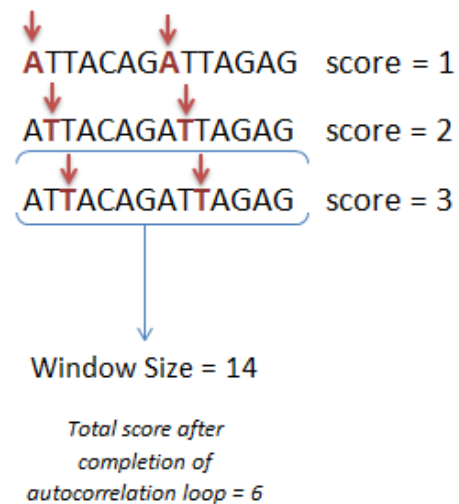*Total score after completion of autocorrelation loop = 6*

Figure 1. Sliding window of autocorrelation function reflecting how scores are incremented within a 7-nt offset and a window size of 14. In this instance, the highest score possible would be 7, so the score of 6 indicates that ~86% of the nucleotides in this window correlated.

nucleotides, it is assigned a score from 1 to 100, and then incremented (and continues with the loop). A score of 1 indicates that there was only 1 matching set of nucleotides within the window (not possible with only four nucleotides), and 100 is a string of the same nucleotide (so that each nucleotide, in return, was correlated with the next). To determine the distribution of scores within any possible genome sequence, a randomized sequence of 100-nt window was run through the autocorrelate function. This process was repeated 100 times using new randomized sequences each time.

Results

The results of the randomized autocorrelation trials allowed for a plot of probabilities, in which the most frequent autocorrelation value was 26.73 with a standard deviation of 3.98 (Figure 2). Of the twenty-four *Bacillus* genomes selected, all of the genomes contained autocorrelation scores over 40, while only nine contained autocorrelation values greater than or equal to 50 (over50 values). Eight of these nine over50 values contained what appeared to be tandem repeats, which ranged from size 4-nt to 11-nt units (mode of 6-nt). Only one of the nine over50 genomes, SPO1, contained multiple instances of autocorrelate values over 50. Five of the over50 values contained tandem repeats were intergenic (the remaining four were outside of gene-coding regions).



Mean = 26.73333
StDev = 3.982144
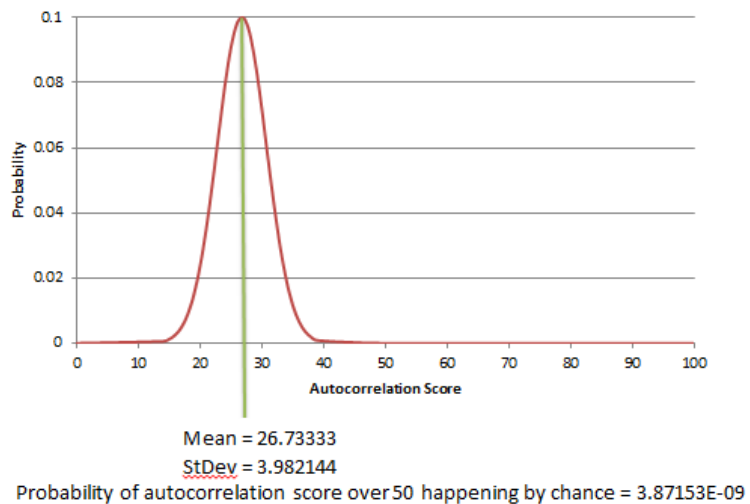Probability of autocorrelation score over 50 happening by chance = 3.87153E-09

Figure 2. Normal distribution of autocorrelation scores. A score of 40 would occur naturally .0390% of the time, while a score of 50 would occur .000000387% of the time.

Discussion

All of the five intergenic over50 values held the purpose of phage replication initiation. Interestingly, this occurrence mirrors the process of *E. coli* replication inititation, in which the initiator protein, dnaA binds to the origin of replication, and thus proceeds to unravel the DNA for replication. What is fascinating however, are the AT-rich tandemly repeated regions that are locations specifically targeted by dnaA[15]. In the *Bacillus* phages studied here, the tandem repeats consisted of 6-nt units; this number differs from the 13-mer repeat found in *E. coli*[15] and 16- and 27-nt repeat units in *Bacillus subtilis*[16]. A reason for a smaller unit size in the observed repeats may be the overall length of the phage genome versus its bacterial genome's counterpart. Repeats in this instance may be proportional to their host's genome size. Another explanation for the decreased unit size amongst *all* of the tandem repeat over50 values may lie in the amount of redundancy needed between genes, where AT-rich regions may serve as attachment sites for infection[13].

Natural occurrences of these repeats would occur by chance less than .000000387% of the time, and thus must have some sort of purpose. Intramolecular recombination may be the cause of these repeats,

where repeats indicate the attP sites needed for host infection. In the remarkable cases described in blue (Appendix Figure 3), where nucleotides seem to be inserted between tandem repeats, literature suggests a type of "cleaning up" conducted by phage genomes. Previous studies have found that, in cases where DNA has been deleted or duplicated by instances of replication errors or slipped-strand mispairing (occurs often in tandem repeats), mutants revert back to wild type phenotypes and proceed to function normally[17]. This may explain the observed blue instances; rather than repeats appearing clearly within the context of the surrounding code, the nucleotides underwent genetic replication errors and failed to maintain tandem repeat status.

Further research in this field may investigate for tandem repeats in not just Bacillus genomes, but all other phage genomes as well. Moreover, focus should be spent on repeats found outside of gene sequences, which cannot currently be explained elaborately by this study.

Determination of Tandem Repeats in Bacillus Bacteriophages
Ellen Korcovelos

Appendix

Figure 3

A51

(28553 50) (28554 50) (28555 50) (28556 50)
(28557 51) (28558 52) (28559 52) (28560 52)
(28561 52) (28562 52) (28563 53) (28564 54)
(28565 54) (28566 54) (28567 53) (28568 52)
(28569 52) (28570 52) (28571 53) (28572 52)
(28573 53) (28574 52) (28575 51) (28576 51)
(28577 50) (28579 50) (28580 50) (28581 50)
(28582 51) (28583 51) (28584 51) (28585 50)
(28586 50) (28587 50) (28588 50) (28589 50)
(28590 50)

Sequence: 28555-28690



Context of:

```
     28555
          |--97--|--159-|
>>BC1888>>-------*------->>BC1889>>    *Between two genes (non-intergentic)

CTT
AAAAAAA
TATACCT
TATA        *Repeats that seem to grow by one nucleotide each time.
TATAA
TATAAA
TATAAAA
C
TAATAAT     *Appears to be a tandem repeat.
TATATAT
TATATAT
ATAGTAT     *Weaker(?) tandem repeats (with above repeat pattern).
TATTTGT
TAATAGT
TATATAG
TATGTAA
TACTATTAATGGTTTAAGGTGTTTGTTTTAGGGGTTTGAT
```
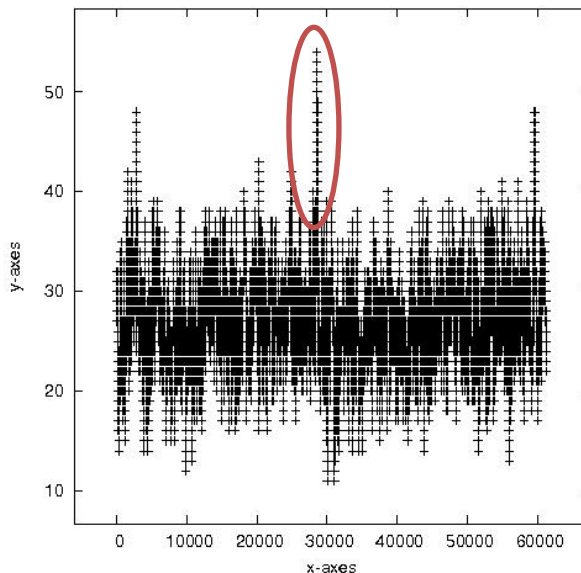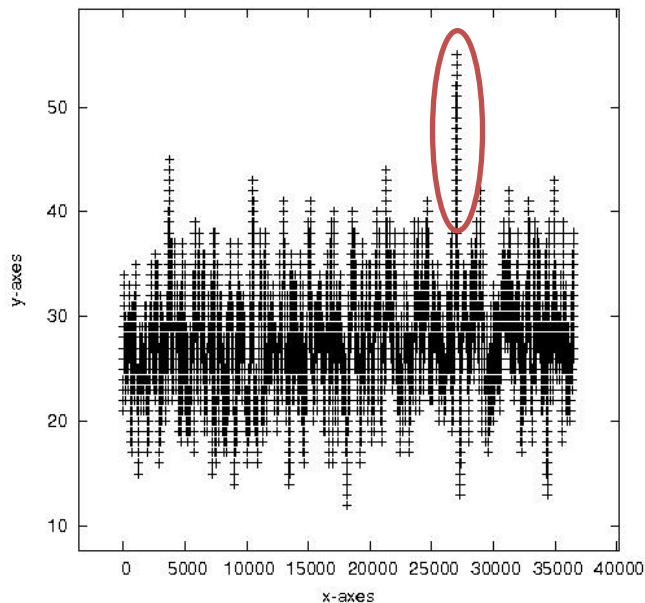
Determination of Tandem Repeats in Bacillus Bacteriophages
Ellen Korcovelos

Cherry

(26977 50) (26978 50) (26983 50) (26984 50)
(26985 50) (26986 50) (26987 51) (26988 51)
(26989 51) (26990 52) (26991 52) (26992 52)
(26993 52) (26994 52) (26995 52) (26996 52)
(26997 52) (26998 51) (26999 51) (27000 51)
(27001 51) (27002 51) (27003 51) (27004 51)
(27005 52) (27006 52) (27007 52) (27008 53)
(27009 53) (27010 53) (27011 54) (27012 54)
(27013 53) (27014 54) (27015 54) (27016 54)
(27017 54) (27018 54) (27019 55) (27020 54)
(27021 53) (27022 52) (27023 52) (27024 52)
(27025 52) (27026 53) (27027 52) (27028 52)
(27029 51) (27030 50) (27031 50)



Sequence: 26977-27131

Context of:

```
    26977
    |--450------|-------518-|
-----&gt;&gt;&gt;&gt;&gt;&gt;CHERRY_0030&gt;&gt;&gt;&gt;&gt;-----        *Intergenic
```

Cherry.CHERRY_0030 (26527 -> 27495)
 Phage replication initiation protein

```
ACGTCCCACGATACGTTAGCGATACGTGACC
AAGAAG        *Appears to be a tandem repeat
AAGAAA
AAGAAC
AAAAAA
AAGAAC
AAAAAG
AAGAAC
AAGAAG
AAAAAG
AAAAAG
AAAAAG
AAAAAC
AAAAAG
AAGAAG
AAAAAG
AACCAG
AAGAAG
AAAAAA
CAAGAA
TAAAATCCAA
```

Determination of Tandem Repeats in Bacillus Bacteriophages
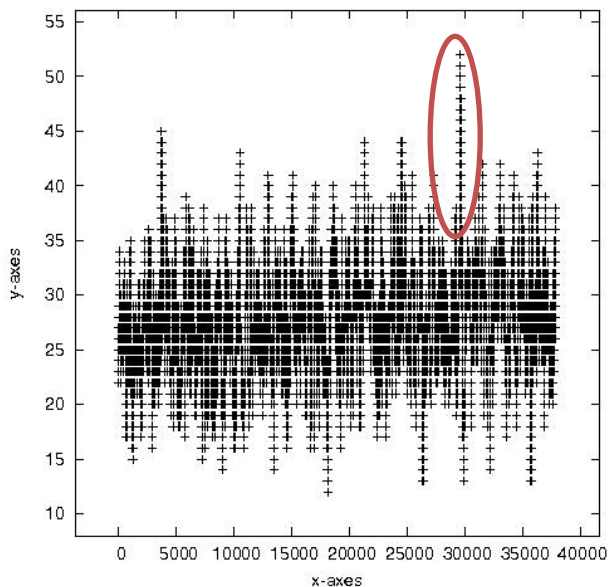Ellen Korcovelos

Fah

(29620 50) (29621 50) (29622 50) (29623 50)
(29624 51) (29625 52) (29626 51) (29627 51)
(29628 51) (29629 50) (29630 51) (29631 51)
(29632 51) (29633 51) (29634 51) (29635 51)
(29636 51) (29637 51) (29638 51) (29639 50)
(29640 50) (29641 50)

Sequence: 29620-29741

Context of:

```
29620
    |--460---|----481-|
----->>>>>>Fah31>>>>>>-----
```



*Intergenic

[Fah.Fah31](#) (29160 -> 30101)
Phage replication initiation protein

```
ATACGTTAGCGATACGTGAC
CAAGAA        *Appears to be a tandem repeat.
GAAGAA
AAAGAA
AAAGAA
AAAGAA
AAAGAA
AAAGAA
AAACAA
AAAGAA
GAAGAA
AAAGAA
CCAGAA        *Weaker(?) tandem repeats (with above repeat pattern).
GAAGAA
AAAACA
AGAATA
AAATCCAAAGCG
```
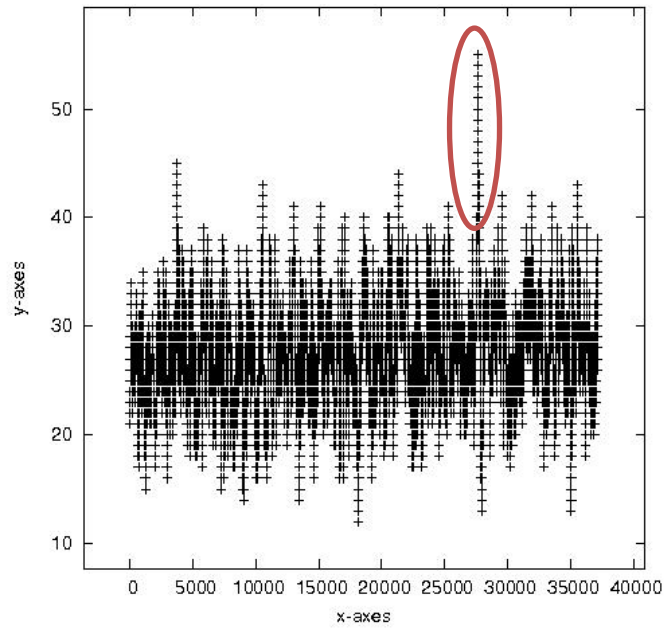
Determination of Tandem Repeats in Bacillus Bacteriophages
Ellen Korcovelos


Gamma51

(27615 50) (27616 50) (27621 50) (27622 50)
(27623 50) (27624 50) (27625 51) (27626 51)
(27627 51) (27628 52) (27629 52) (27630 52)
(27631 52) (27632 52) (27633 52) (27634 52)
(27635 52) (27636 51) (27637 51) (27638 51)
(27639 51) (27640 51) (27641 51) (27642 51)
(27643 52) (27644 52) (27645 52) (27646 53)
(27647 53) (27648 53) (27649 54) (27650 54)
(27651 53) (27652 54) (27653 54) (27654 54)
(27655 54) (27656 54) (27657 55) (27658 54)
(27659 53) (27660 52) (27661 52) (27662 52)
(27663 52) (27664 53) (27665 52) (27666 52)
(27667 51) (27668 50) (27669 50)

Sequence: 27615-27769

Context of:

```
                 27615
    |--450--------|--------518-|
------>>>>>>GAMMAUSAM_0032>>>>>>-----          *Intergenic
```

Bac-gamma-51.GAMMAUSAM_0032 (27165 -> 28133)
 phage replisome organizer protein,
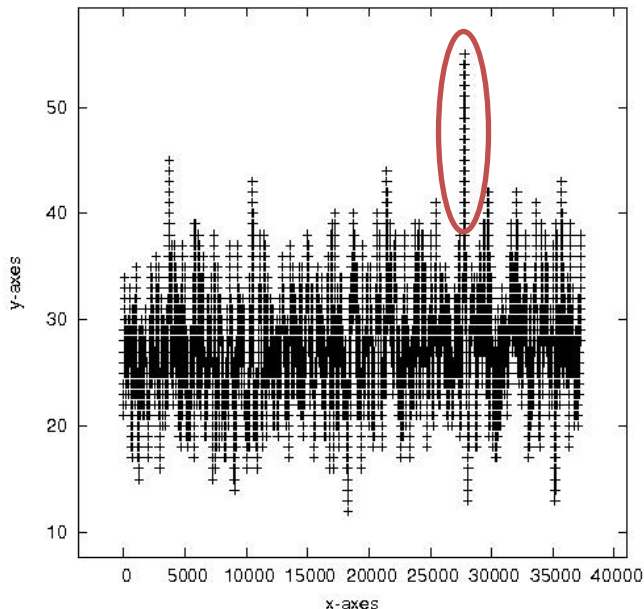 putative

```
ACGTCCCACGATACGTTAGCGATACGTGAC
CAAGAA          *Weaker, but still apparent tandem repeats.
GAAGAA
AAAGAA
CAAAAA
AAAGAA
CAAAAA
GAAGAA
CAAGAA
GAAAAA          *Very strong appearance of tandem repeats.
GAAAAA
GAAAAA
GAAAAA
CAAAAA
GAAGAA
GAAAAA
GAACCA
GAAGAA
GAAAAA
ACAAGAATAAAATCCAA
```

# Determination of Tandem Repeats in Bacillus Bacteriophages
Ellen Korcovelos

## GammaH

(27744 50) (27745 50) (27750 50) (27751 50)
(27752 50) (27753 50) (27754 51) (27755 51)
(27756 51) (27757 52) (27758 52) (27759 52)
(27760 52) (27761 52) (27762 52) (27763 52)
(27764 52) (27765 51) (27766 51) (27767 51)
(27768 51) (27769 51) (27770 51) (27771 51)
(27772 52) (27773 52) (27774 52) (27775 53)
(27776 53) (27777 53) (27778 54) (27779 54)
(27780 53) (27781 54) (27782 54) (27783 54)
(27784 54) (27785 54) (27786 55) (27787 54)
(27788 53) (27789 52) (27790 52) (27791 52)
(27792 52) (27793 53) (27794 52) (27795 52)
(27796 51) (27797 50) (27798 50)



Sequence: 27744-27898
Context of:
```
            27744
   |--429------|-------518-|
----->>>>>>dHerelle-33>>>>>>-----        *Intergenic
```

dHerelle.dHerelle-33 (27315 -> 28262)
 Phage replication initiation protein

```
ACGTCCCACGATACGTTAGCGATACGTGAC
CAAGAA      *Weaker, but still apparent tandem repeats.
GAAGAA
AAAGAA
CAAAAA
AAAGAA
CAAAAA
GAAGAA
CAAGAA
GAAAAA      *Appears to be a highly conserved tandem repeat.
GAAAAA
GAAAAA
GAAAAA
CAAAAA
GAAGAA
GAAAAA
GAACCA
GAAGAA
GAAAAA
ACAAGAATAAAATCCAA
```

# Determination of Tandem Repeats in Bacillus Bacteriophages
Ellen Korcovelos

Spo1
(11332 51) (11333 52) (11334 52) (11335 52)
(11336 53) (11337 52) (11338 51) (11342 51)
(11343 51) (44112 51) (44113 51) (44123 51)
(44124 52) (44125 51) (78296 51) (78297 51)
(78300 51) (78301 51) (78318 51) (78320 51)
(78321 52) (78322 52) (78323 52) (78324 52)
(78325 53) (78326 53) (78327 53) (78328 53)
(78329 53) (78330 53) (78331 53) (78332 53)
(78333 53) (78334 54) (78335 54) (78336 53)
(78337 53) (78338 54) (78339 53) (78340 52)
(78341 52) (78342 53) (78343 53) (78344 53)
(78345 53) (78346 52) (78347 51) (78348 51)
(78349 51) (78350 51) (78351 51))



Sequence: 11332-11443
Context of:

```
                11332
          |--183-|--205-|
<<SPO1_25<<-------*-------<<SPO1_26<<          *Non-intergenic
```

AAGT                    *Again, these sequences appear to increment by one adenine
AAAGAG                  in the repeat (however, this one does not appear to be
AAAAATT                 tandem)
AAAAGA
TTTTTTAAGTC             *One series of possible tandem repeats.
TTTTTAAAAGC
TTTTTTAAAGT
ACTTTAAAAGA
TTTTTTTTTTT             *Highly conserved tandem repeat.
TTTTTTTTTTT
TTTTTTTTTTT
TTTTTTTTTTT
A


Sequence: 44112-44225
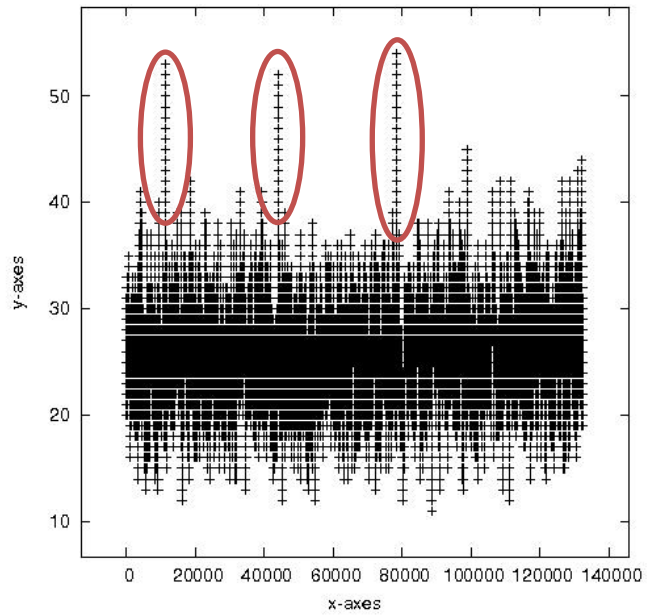Content of:

```
                44112
          |--15--|--90--|
>>SPO1_77>>-------*------->>SPO1_78>>          *Intergenic
```

AGATAGGAGACAGCTAAGTGCTGTCTCCTA
TTTT
TGTT
CCCT
TCTG
TTTTCCAGA               *Very unlikely to be tandem repeat (but may explain
TTTTAGATT                correlation value)
ATAGAC TATAAA [promoter sequence]GGAGGTAATGGGAAATGCCAAAAATCCAGAATGTAGCA
[start of gene 78]

# Determination of Tandem Repeats in Bacillus Bacteriophages
Ellen Korcovelos


Sequence: 78296-78451
Context of:

```
                78296
           |--35--|--183-|
>>SPO1_109>>-------*------->>SPO1_110>>        *Non-intergenic
```

```
TAGTATA
AAAAGTGAAAA
AGAGACTTTT          *Appears to be a tandem repeat, however nucleotides are
AAAGACCTTTT           added
TAAAAGCTTTAA
AAAAAAACTTTAAA
AGAA
AAAG
AAAA               *Highly conserved tandem repeats
AAAA
AAAA
AAAA
AAAA
AAAA
AAAA
AAAA
AATT               *Maybe be parts of the tandem repeats
TATT
TTTG
TCCTACGAAGTAGGACAAAAAGAGCGCAACACAACAATA
```

Determination of Tandem Repeats in Bacillus Bacteriophages
Ellen Korcovelos

IEBH

(16464 50) (16465 50) (16466 50) (16467 50)
(16468 50) (16469 51) (16470 50) (16471 51)
(16472 50) (16473 50)



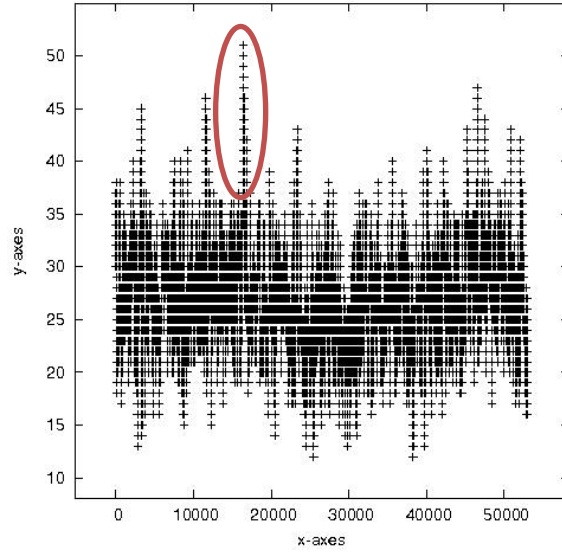Sequence: 16464-16573
```
Context of:
                16464
             |--905-|--142-|
>>IEBH_gp41>>-------*------->>IEBH_gp42>>      *Non-intergenic

TGATGGC                            *There does not appear to be any
TTTTTG                              striking repeats here.
TTTTTTGTAAACTAGGG
ATTATCA
ATTATGTTAGTTCCTAGTTTAGAGAGAATAAA
ATTATTTACTATATAGAA
ATTACACATTAAACGTAAAACAG
```

Determination of Tandem Repeats in Bacillus Bacteriophages
Ellen Korcovelos

WBeta

(29776 50) (29777 50) (29782 50) (29783 50)
(29784 50) (29785 50) (29786 51) (29787 51)
(29788 51) (29789 52) (29790 52) (29791 52)
(29792 52) (29793 52) (29794 52) (29795 52)
(29796 52) (29797 51) (29798 51) (29799 51)
(29800 51) (29801 51) (29802 51) (29803 51)
(29804 52) (29805 52) (29806 52) (29807 53)
(29808 53) (29809 53) (29810 54) (29811 54)
(29812 53) (29813 54) (29814 54) (29815 54)
(29816 54) (29817 54) (29818 55) (29819 54)
(29820 53) (29821 52) (29822 52) (29823 52)
(29824 52) (29825 53) (29826 52) (29827 52)
(29828 51) (29829 50) (29830 50)



Sequence: 29776-29830
Context of:
```
            29776
    |--429-----|-----518-|
----->>>>>>WBeta-34>>>>>>-----          *Intergenic
```
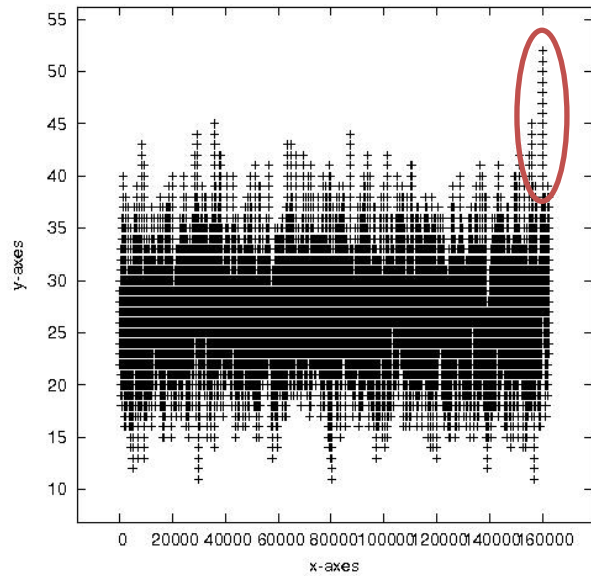
WBeta.WBeta-34 (29347 -> 30294)

ACGTCCCACGAT        *Seems to be a CRISPR.
ACGTTAGCGAT
ACGTGACC
AAGAAG              *Appears to be a tandem repeat.
AAGAAA
AAGAAC
AAAAAA

## B4

(160019 50) (160020 50) (160021 51) (160022 50)
(160026 50) (160027 50) (160028 51) (160029 52)
(160030 52) (160031 51) (160032 51) (160033 51)
(160034 50) (160035 50) (160036 51) (160037 50)
(160038 51) (160039 51) (160040 51) (160041 51)
(160042 51) (160043 50) (160048 50) (160049 51)
(160050 51) (160051 50) (160052 50) (160053 50)
(160054 50) (160055 50) (160056 50) (160057 50)
(160068 50) (160069 50) (160070 51) (160071 52)
(160072 51) (160073 51) (160074 51) (160075 50)
(160076 50) (160077 50) (160078 50) (160094 50)
(160095 50) (160103 50)



Sequence: 160019-160203
Context of:

```
                  160019
             |--557-|--428-|
<<BacB4-0278<<-------*------->>BacB4-0279>>    *Non-intergenic
```
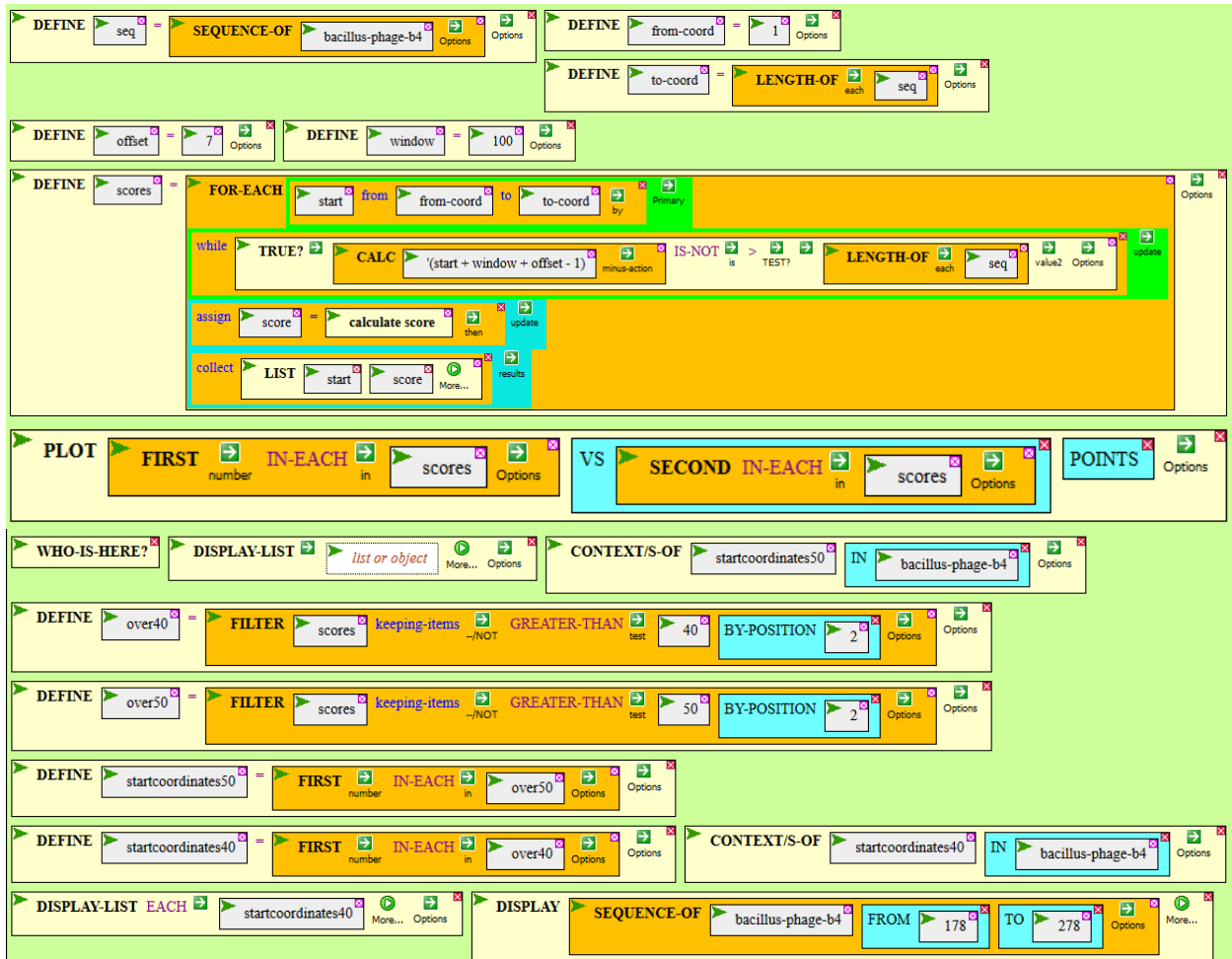
```
G
TATAATAAA         *Appears to be a tandem repeat.
TATAACTAA
ATAAAATAA
CTAATATTA
AATATTATA         *Appears to be a weaker repeat of the sequence above.
TATGTAATT
AAATATAAA
TAATAAATA
CACTAATAGATATAGGGAA
CAGTTAAAAATAAAATAGGGAATAGTATCAGGAAT
ATTACTT           *Appears to be a weak repeat.
ATTATTT
ATTTAAA
AATTAAT
ACTTTAG
TCAGAATAAGAATTAAGGATATG
```

Determination of Tandem Repeats in Bacillus Bacteriophages
Ellen Korcovelos

Figure 4

References

1. Scherzinger et al., Huntingtin-encoded polyglutamine expansions form amyloid-like protein aggregates in vitro and in vivo. *Cell.* 1997, 90(3):549-558.
2. Lee JE, Cooper TA. Pathogenic mechanisms of myotonic dystrophy. *Biochem Soc Trans.* 2009, 37:1281-1286.
3. Muragaki Y, Mundlos S, Upton J, Olsen BR. Altered growth and branching patterns in synpolydactlyy caused by mutations in HOXD13. *Science.* 1996, 272(5261):548-551.
4. Lindstedt B, Heir E, Gjernes E, Vardund T, Kapperud G. DNA Fingerprinting of Shiga-toxin producing Escherichia coli O157 based on multiple-locus variable-number tandem-repeats analysis (MLVA). *Annals of Clinical Microbiology and Antimicrobials.* 2003, 2:12.
5. Mazel D, Houmard J, Castets AM, and Tandeu N. Highly repetitive DNA sequences in cyanobacterial genomes. *Journal of Bacteriology.* 1990, 172(5):2755-2761.
6. Lupski J, Weinstock G. Short, interspersed repetitive DNA sequences in prokaryotic genomes. *Journal of Bacteriology.* July, 174(14):4525-4529.
7. Van Belkum A, Scherer S, van Alphen L, Verbrugh H. Short-sequence DNA repeats in prokaryotic genomes. *Microbiol Mol Biol Review.* 1998, 6(2):275-293.
8. Hillyar C. Genetic recombination in bacteriophage lambda. *Bioscience Horizons.* 2012, 5.
9. Okubo S, Strauss B, Stodolsky M. The possible role of recombination in the infection of competent Bacillus subtilis by bacteriophage deoxyribonucleic acid. *Virology.* 1964, 24(4):552-562.
10. Morris P, Marinelli L, Jacobs-Sera D, Hendrix R, Hatfull G. Genomic characterization of mycobacteriophage Giles: evidence for phage acquisition of host DNA by illegitimate recombination. *Journal of Bacteriology.* 2008, 190(6):2172.
11. Pope et al. Expanding the diversity of mycobacteriophages: insights into the genome architecture and evolution. *PLoS ONE.* 2012, 6(1):e16329.
12. Mesnage S, Weber-Levy M, Haustant M, Mock M, Fouet A. Cell surface-exposed tetanus toxin fragment C produced by recombinant Bacillus anthracis protects against tetanus toxin. *Infection and Immunity.* 1999, 67(9):4847-4850.
13. Alberts B, et al. Site-Specific Recombination. *Molecular Biology of the Cell.* 4th edition. New York: Garland Science. 2002.
14. Program written by Jeffrey Elhai. http://biobike.csbc.vcu.edu/
15. Bramhill D, Kornberg A. A model for initiation at origins of DNA replication. *Cell.* 1988, 54(7):915-918.
16. Krause M, Ruckert B, Lurz R, Messer W. Complexes at the replication origin of Bacillus subtilis with homologous and heterologous dnaA protein. *Journal of Mol Bio.* 1997, 274(3):365-380.
17. Sastalla I, Rosovitz M, Leppla S. Accidental selection and intentional restoration of sporulation-deficient anthracis mutants. *Applied and Environmental Microbio.* 2010, 76(18):6318-6321.