

A tandem repeats database for bacterial genomes: application to the genotyping of *Yersinia pestis* and *Bacillus anthracis*

ABSTRACT

Bacillus anthracis and *Yersinia pestis* are classified as having an average density of tandem repeat arrays over 100 bp in comparison to other bacterial genomes studied so far, and testing a fraction of these sequences for polymorphism yielded dozens of informative markers, some of which displayed remarkably high levels of polymorphism, including allele length from 600-1950 bps and nine alleles resolved in the small number of independent *Bacillus* strains typed here.

INTRODUCTION

A tandem repeats database for bacterial genomes (TDB) is described. The TDB contains the sequences of all the tandem repeat genes of the *E. coli* strain isolates of the genus *Bacillus anthracis* (*B. anthracis*), and of the human isolates of the genus *Pseudomonas aeruginosa* (*Pseudomonas aeruginosa*), and of the *E. coli* strain isolates of the genus *Streptococcus thermophilus* and the human isolates of the genus *Listeria monocytogenes*. The TDB also contains the human isolates of the genus *Listeria monocytogenes* and *Streptococcus thermophilus* from a previous study by *B. anthracis* isolates. Tandem repeats are typically classified into satellites (encompassing megabases of DNA, associated with heterochromatin) as well as minisatellites (repeat units in the range of 6-100 bp, spanning hundreds of base-pairs) and microsatellites (reasonably placed by a few tens to dozens of nucleotides), making them an essential component of mammalian genetics for creating genetic maps. A number of studies have supported the idea that tandem repeats, such as mini and microsatellites, can serve as a crucial source of informative markers for identification of pathogenic bacteria. This suggests that these markers play varying roles in phenotypic variation within bacteria, with some occurring within the regulatory region or within coding regions where the unit length is not more than three times longer (reviewed in). Human and yeast eukaryotes have been studied for their mutating processes towards micro and minisatellites, with the results suggesting that these mechanisms are more complex than those of other organisms. In bacteria, the simple sequence contingency loci, which are those with a tandem repeat from the microsatellite class (repeat unit sizes of 1-8 bp), is useful for achieving reversible on and off states of expression for the corresponding gene. The mutation rate of tetranucleotide (*Haemophilus influenzae*) is above 10^{-4} , and different alleles are observed in bacterial strains. Minisatellite's array is also tested for identification and phys The report will initially describe the utilization of a tandem repeats database for bacterial genomes and briefly compare the general features of tandem repetitions in various bacterium genome-cloned strains with publicly available sequences. Furthermore, it will demonstrate how this tool can be applied to rapidly characterizing highly polymorphic markers in two pathogens, *Y. pestis* and *B. anthracis*. The ancestors of both *Y. pestis* and *B. anthracis*, which are classified as causative agents of plague and anthrax, have recently emerged. However, RFLP analysis of IS100 locations is more difficult than PCR typing, making it the most suitable technology for typing a high-resolution yeast species. In the case of *Bryx siciliana*, early polymorphisms were identified essentially using AFLP (Amplified Fragment Length Polymorphic

CONCLUSION

Remarkable conclusions The focus of this investigation was on minisatellites, which are repeat units longer than 9 base-pairs, so as to avoid duplication of work and nomenclature, and to facilitate the typing alleles with agarose gel electrophoresis. However, simple sequence contingency loci are also represented in the database and are of great interest for molecular pathogenicity studies. The tandem repeats database was used here for identification and epidemiological analyses among human pathogen species that are amenable to these types of humans. The density and frequency of tandem repeats differ among bacterial species, with some having an overabundance of these repeat units with multiple of three-units. Figure 7 demonstrates that the two species studied exhibit high levels of polymorphism in tandem repeats, which are strongly correlated with sequenced allele characteristics. However, there is no clear correlation between these traits and the overall homogeneity of the tandem array in *Bacillus anthracis*, and even more significant differences between microsatellites (1-3 bp repeat units) may prevent this phenomenon from being explained through systematic allele sequencing. *Cicadaminers* belonging to the *B. anthracis* marker class (Ceb-Bams1, 3, 7, 13, and 30) are highly polymorphic, with PIC values above 0.7, and thus the allele length measured for Ceb-Bams1 in the Ames strain is not of that magnitude as expected from the sequence data (Table 2). This result may be due to a high mutation rate at Cebit-Bummer, or corresponding sequencing error: the estimated allele size corresponds to nothingness, but only three alleles. Despite the lack of consideration for allele size differences in the distance matrix, the phylogenetic tree depicted in Figure 6 often clusters strains with alleles of similar size at the most variable loci. This suggests that mutation events are primarily small-scale changes and necessitates more comprehensive studies using full allele sequencing to better understand the sequence of events that result in a population of alleles.