DNA BASE-CALLING TECHNIQUES


Fadi Odeh, Candidate for the Master of Science Degree


American University of Sharjah, 2008


ABSTRACT

The availability of substantial amounts of DNA sequence information has begun to revolutionize the practice of biology. So it is obvious that manual sequencing output is not adequate to keep pace with the growing demand and is far from what is required to obtain the 3-billion-base human genome sequence.

To avoid this difficulty, replacing manual sequencing with an automated one is essential, and it is particularly important that human involvement in data processing be significantly reduced or eliminated. Progress in this respect requires both improving the amount of error-free data being processed, as well as the reliable accuracy measures to reduce the need for human involvement in error correction. Here, we precede one step toward that goal: a basecalling program for automated DNA sequencing, with improved accuracy.

The major goal of this thesis is to develop a new basecalling technique to improve the efficiency of the DNA sequencing process. Improved efficiency will be achieved by increasing the average length of error-free sequences and enhancing the base identification

process at the beginning and end of the DNA sequences. This will greatly increase sequencing throughput and reduce both cost and error associated with the current DNA sequencing process. ABI machines (Applied Bio-systems Incorporated sequencing machines) are currently the major source of reading DNA data. They are capable of producing sequences of 1000 bases in length (bases produced by PCR (Polymerase chain reaction)). These machines are associated with basecalling software, the most advanced software is called KB Basecaller v1.4 and it is publicly used by the sequencing community because of its reliability and accuracy. It can produce impressive results of 500~600 error-free sequences. The error-free sequences are normally located in the middle of the 1000 base length where the data is clear, and bases are easily distinguishable. However, the bases at the beginning and end of a 1000 base sequence are obscure and difficult to identify. The base calling error in these regions is relatively high. Thus the average basecalling error over a 1000 base sequence is between 3.5 and 6%. The foundation of this proposed research is based on a new base-calling program related to combining signal processing and pattern recognition systems which includes the following steps: noise filtration, baseline adjustment, mobility shift correction, feature extraction and the development of an intelligent basecalling algorithm. The new algorithm will be tested and validated on a number of pre-sequenced DNA sequences.

Combining Gaussian Mixture Models and Hidden Markov Models (GMM-HMM) classifier will be used as a classification model for the recognition of the DNA bases based on its several advantages over other classifiers in that they do not require heavy training, they are very simple to implement with the number of classes, and they ensure the coverage of the statistical properties of the data using Gaussian distribution.

DNA sequence information is critical to understand genetic variations that can influence both disease, and genetic interactions, which in turn can influence drug efficacy. As such, automated sequencers play a vital role in the drug discovery process.

CONTENTS

LIST OF FIGURES