# Assignment 1: Abstract and Introduction Summary

## Abstract :

Tiny non-coding RNAs generated by eukaryotes serve as specificity determinants for various gene-regulatory complexes. MicroRNAs (miRNAs), endogenous short interfering RNAs (siRNAs), and Piwi-associated RNAs are examples of these (piRNAs). Small RNA-seq, a variant RNA-seq approach based on highly parallel sequencing, was used to annotate and quantify the data. The alignment of small RNA-seq data to a reference genome is a crucial step in the research process. Due to their small size (20-30 nts depending on the organism and sub-type) and proclivity to stick together, Reads that align equally well to multiple genomic locations are very common, particularly if they come from multi-gene families or repeated regions. Typical approaches to dealing with multi-mapped narrow RNA-seq reads trade off precision for sensitivity. The method 'butter' combines accuracy and sensitivity by using an iterative approach to place multi-mapped reads, where the choice between possible positions is made each time. Dictated by the densities of more confidently matched reads in the surrounding area. In comparison to other small RNA-seq aligners, Butter performs well.

## Introduction :

Small RNA-seq is a type of RNA-seq that uses directional RNA ligation, cDNA synthesis, and sequencing to capture small RNAs, usually in the 15-40 nt range. It's an effective tool for finding, annotating, and quantifying small RNA molecules. Multi-mapped reads, known as reads with more than one equally probable alignment location, are a common problem in RNA-seq data alignment, particularly for small data sets. Some of solution is to (completely disregard multi-mapped reads, place multi-mapped reads randomly at one of the possible positions. simply report all possible alignment positions).

Larger RNA precursors are converted into piRNAs, endogenous siRNAs, and miRNAs. Processing is heterogeneous to differing degrees, resulting in several

distinct small RNAs that, when properly mapped to the reference genome, provide a cluster of alignments at the origin locus.

I present butter (Bowtie Using iterative placement of Repetitive small rnas), a technique that employs an iterative approach based on uniquely and low-level multi-mapped reads to improve the precision of small RNA-seq genomic alignments.

## Related works:

Multi-mapped reads, described as reads with more than one equally likely alignment location, are a common problem in RNA-seq data alignment, particularly for small data sets. One easy solution is to completely disregard multi-mapped reads. For canonical RNA-seq data (e.g. fragments of polyA+ mRNAs), this method is commonly used. The number of multi-mapped reads in canonical RNA-seq data is minimized due to relatively long read lengths (150 nts), paired-end sequencing, and the fact that most polyA+ mRNAs are inherently non-repetitive. None of these caveats relate to limited RNA-seq, and missing multi-mapped reads completely results in a significant loss of data. Another straightforward approach is to randomly insert multi-mapped reads in one of the available locations. Many common aligners use this behavior by default. However, this approach has the apparent disadvantage of causing most multi-mapped reads to be incorrectly placed. A read with two possible locations, for example, has a 50% chance of being placed incorrectly; a read with three possible locations has a 67 percent chance of being placed incorrectly, and so on. A third straightforward solution is to simply list all possible alignment positions. However, this results in the alignment files becoming excessively large. Since the alignments are no longer indicative of abundance, this approach is vulnerable to misinterpretation in downstream studies, resulting in overestimation of small RNA output from repetitive areas of the genome.

## Results and methodology:

**Butter outperforms all other tested alignment methods in balancing sensitivity with precision**

Using five sets of simulated short RNA-seq data from Arabidopsis, rice, and maize, the performance of butter was compared to seven alternative alignment techniques. Because, unlike actual data, the proper alignment site for each simulated read was known with confidence, simulated data were employed in this research. As a consequence, each read's outcome may be classified as a true positive (TP; read aligned to the proper place), a false positive (FP; read aligned to the erroneous site), or a false negative (FN; read aligned to the erroneous position) (FN; a read that was not aligned to the genome). There were no true negatives in the simulated tiny RNA-seq data since there were no reads that should not have aligned. For multi-mapped reads, the earlier version of ShortStack (version 1.2.4; (16), NoVo align, bwa (17), and bowtie (12) all choose and report just one alignment location at random. For multi-mapped reads, Patman (18) and cashx searchDB (19) report all potential alignments. Butter selects a single alignment location for multimapped reads based on the relative densities of unique and confidently placed reads, as well as suppressing the reporting of any location for multimapped reads for which a high-confidence choice cannot be reached (Figure 1; See methods). Finally, bowtie suppresses all alignments for multi-mapped reads when the -m option is set to 1.

**Performance of butter with real data**

Using actual data sets from Arabidopsis inflorescence tissue (sequence read archive accession SRR1042171; 13), rice lamina joints from flag leaves (SRR976171; 14), and maize leaves (SRR1186264; 15), the influence of alignment approach on small RNA-seq data interpretation was investigated. These datasets have depths of 14.3, 31.1, and 85.7 million adapter-trimmed reads, respectively, and are indicative of contemporary small RNA-seq research. Bowtie was compared to butter using two frequently used alignment methods: bowtie with default settings, which reports just one location for multi-mapped reads at random, and bowtie with option -m set to 1, which suppresses all alignments for multi-mapping reads. Arabidopsis had the fewest multi-mapped reads, whereas maize had the most, as predicted based on genomic architectures and sizes. Surprisingly, clustering density could be used to locate virtually all multi-mapped reads in all three species.