

Sequence analysis

PRAPI: post-transcriptional regulation analysis pipeline for Iso-Seq

Yubang Gao¹, Huiyuan Wang¹, Hangxiao Zhang¹, Yongsheng Wang¹,
Jinfeng Chen² and Lianfeng Gu^{1,*}

¹Basic Forestry and Proteomics Research Center, Fujian Provincial Key Laboratory of Haixia Applied Plant Systems Biology, College of Life Science, Fujian Agriculture and Forestry University, Fuzhou 350002, China and
²Department of Plant Pathology and Microbiology, Institute of Integrative Genome Biology, University of California, Riverside, CA 92521, USA

*To whom correspondence should be addressed.
Associate Editor: John Hancock

Received on September 4, 2017; revised on October 24, 2017; editorial decision on December 16, 2017; accepted on December 20, 2017

Abstract

Summary: The single-molecule real-time (SMRT) isoform sequencing (Iso-Seq) based on Pacific Bioscience (PacBio) platform has received increasing attention for its ability to explore full-length isoforms. Thus, comprehensive tools for Iso-Seq bioinformatics analysis are extremely useful. Here, we present a one-stop solution for Iso-Seq analysis, called PRAPI to analyze alternative transcription initiation (ATI), alternative splicing (AS), alternative cleavage and polyadenylation (APA), natural antisense transcripts (NAT), and circular RNAs (circRNAs) comprehensively. PRAPI is capable of combining Iso-Seq full-length isoforms with short read data, such as RNA-Seq or polyadenylation site sequencing (PAS-seq) for differential expression analysis of NAT, AS, APA and circRNAs. Furthermore, PRAPI can annotate new genes and correct mis-annotated genes when gene annotation is available. Finally, PRAPI generates high-quality vector graphics to visualize and highlight the Iso-Seq results.

Availability and implementation: The Dockerfile of PRAPI is available at <http://www.bioinfor.org/tool/PRAPI>.

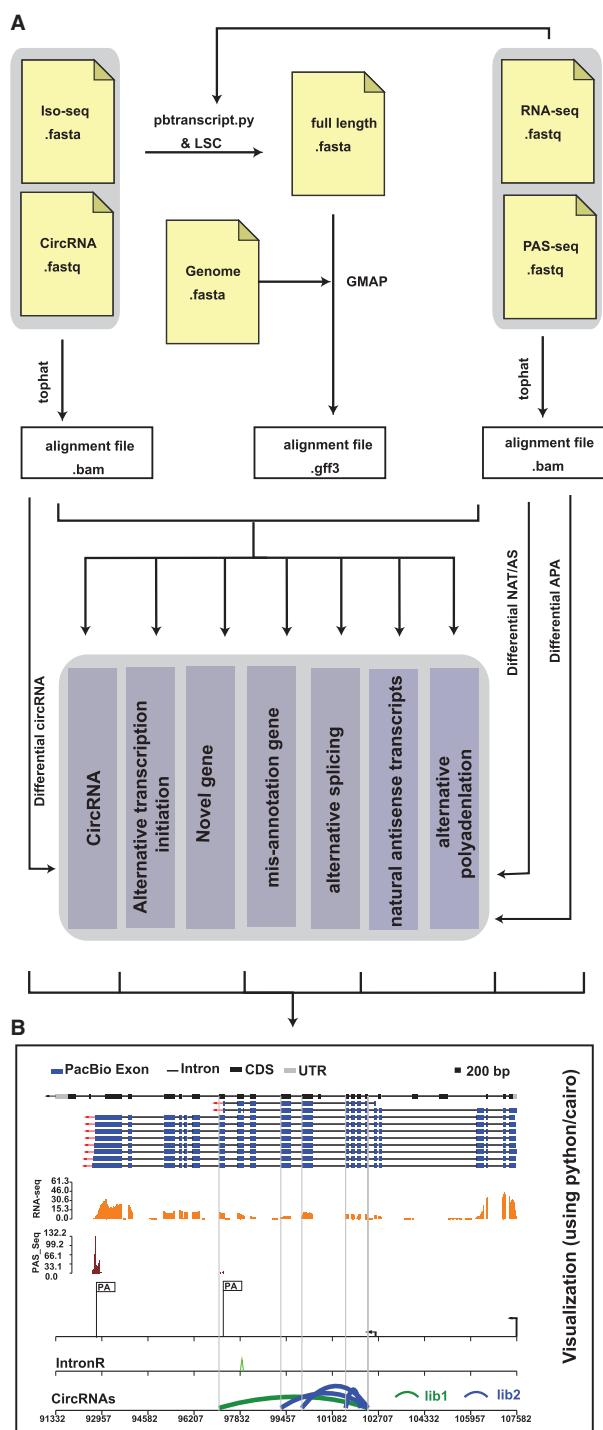
Contact: lfgu@fafu.edu.cn

1 Introduction

PacBio platform offers a reliable way to identify alternative splicing (AS) and alternative cleavage and polyadenylation (APA) from Iso-Seq (Abdel-Ghany *et al.*, 2016; Wang *et al.*, 2017). To the best of our knowledge, TAPIS (Abdel-Ghany *et al.*, 2016), IDP (Au *et al.*, 2013), and PacBio's SMRT-Analysis are excellent tools for analyzing isoforms from Iso-Seq reads. TAPIS also performs APA analysis since full-length isoforms of Iso-Seq span from the transcription start sites to terminal sites (Abdel-Ghany *et al.*, 2016). At present, transcriptome profiling with next-generation sequencing (NGS) data such as RNA-Seq and polyadenylation site sequencing (PAS-Seq), is an indispensable method for performing quantitative analysis of gene or isoform expression (Wang *et al.*, 2017). However, previous

methods lack the module to combine high-accuracy short reads with Iso-Seq long reads. Consequently, a full pipeline consisting of Iso-Seq data process, comprehensive analysis, and visualization, will be necessary when the cost of Iso-Seq is reduced. To meet the requirement of Iso-Seq analysis, we have developed a pipeline called Post-transcriptional Regulation Analysis Pipeline for Iso-Seq (PRAPI), which aims to identify and quantify the post-transcriptional regulation. PRAPI also adds new features to allow users to identify differentially expressed NAT, AS, APA and circRNAs by further combining them with short reads sequencing. Furthermore, the results from PRAPI are highlighted using vector diagrams. In summary, PRAPI is a comprehensive, user-friendly software that facilitates Iso-Seq analysis.

Time consumed for PRAPI was evaluated using the data from a previous study (Wang *et al.*, 2017), which includes 146 225 error-



corrected full-length non-chimeric reads. The time required for processing the whole dataset is 98 min using 18.0 GB Memory on our server (Linux: Intel(R) Xeon(R) CPU E5-2660 v3 @ 2.60 GHz; 20 Cores; Memory 256GB).

In summary, we present PRAPI, a full-function pipeline aimed at not only interpreting post-transcriptional regulation based on PacBio's

full-length reads but also presenting the visualization of the high throughput data. Descriptions of all the input, output, and parameters can be found in our online tutorials and test dataset. The source code with complete tutorial is publicly available.

Acknowledgement

The authors thank Professor Anireddy S.N. Reddy and Dr. Michael Hamilton for advice on TAPIS.

Funding

This work was supported by the National Key Research and Development Program of China (2016YFD0600106), the National Natural Science Foundation of China Grant (31570674).

Conflict of Interest: none declared.

References

Abdel-Ghany, S.E. *et al.* (2016) A survey of the sorghum transcriptome using single-molecule long reads. *Nat. Commun.*, **7**, 11706.

- Au, K.F. *et al.* (2013) Characterization of the human ESC transcriptome by hybrid sequencing. *Proc. Nat. Acad. Sci. USA.*, **110**, E4821–E4830.
- Au, K.F. *et al.* (2012) Improving PacBio long read accuracy by short read alignment. *PLoS One*, **7**, e46679.
- Gao, Y. *et al.* (2015) CIRI: an efficient and unbiased algorithm for de novo circular RNA identification. *Genome Biol.*, **16**, 4.
- Li, H. *et al.* (2009) The sequence alignment/map format and SAMtools. *Bioinformatics*, **25**, 2078–2079.
- Robinson, M.D. *et al.* (2010) edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics*, **26**, 139–140.
- Shen, S. *et al.* (2014) rMATS: robust and flexible detection of differential alternative splicing from replicate RNA-Seq data. *Proc Natl Acad. Sci. USA.*, **111**, E5593–E5601.
- Wang, L. *et al.* (2010) DEGseq: an R package for identifying differentially expressed genes from RNA-seq data. *Bioinformatics*, **26**, 136–138.
- Wang, T. *et al.* (2017) Comprehensive profiling of rhizome-associated alternative splicing and alternative polyadenylation in moso bamboo (*Phyllostachys edulis*). *Plant J.*, **91**, 684–699.
- Wu, T.D., and Watanabe, C.K. (2005) GMAP: a genomic mapping and alignment program for mRNA and EST sequences. *Bioinformatics*, **21**, 1859–1875.