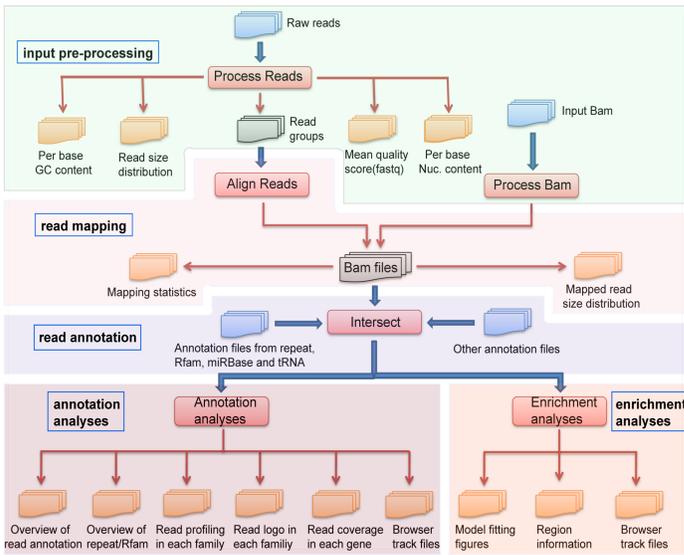


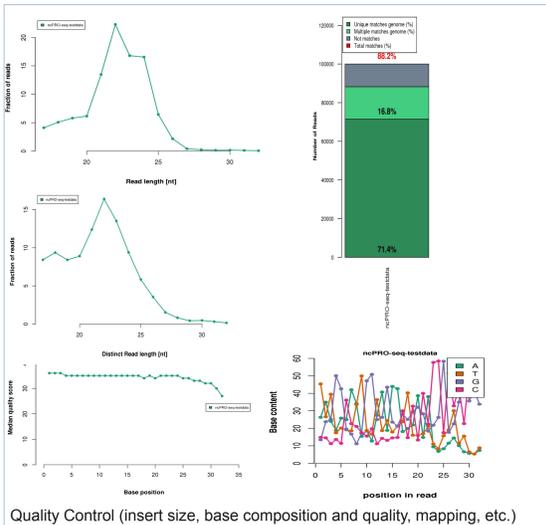
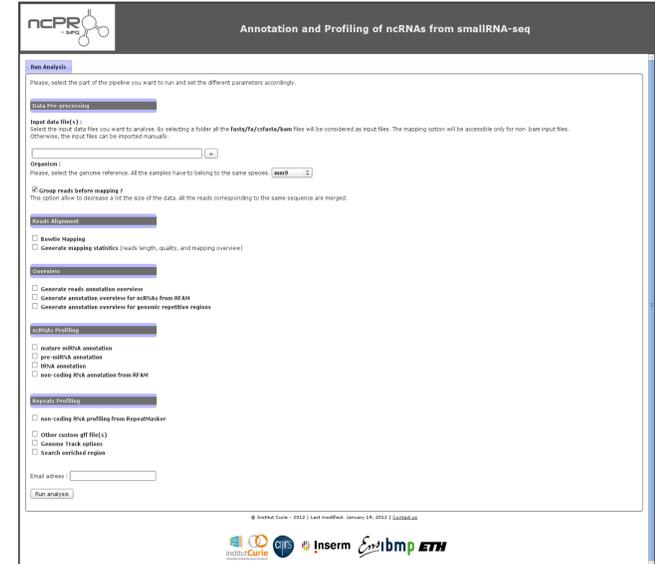
Over recent years, deep sequencing technology has become a powerful approach for investigating **small non-coding RNA** (ncRNA) populations, i.e. small RNA-seq. It is now established that an increasing number of novel small ncRNA families distinct from microRNAs are generated over kingdoms from different coding/non coding regions via various biogenesis pathways and might involve a great spectrum of biological processes. For example, two other major classes of endogenous small RNAs, Piwi-interacting RNAs (piRNAs) and endogenous small interfering RNAs (endo-siRNAs), have been identified and widely investigated in mammals [1]. Moreover, in other organisms like plants more classes of small ncRNA have been described indicating that a wide range of small ncRNAs exist [2]. However, most of the existing tools devoted to sRNA-seq analysis, are only based on miRNAs annotation and quantification, significantly neglecting other types or new types of small ncRNAs. Moreover, they just perform gene-based analysis, but not detailed family-based (profiling) analysis which is critically important to investigate known small ncRNA families and to identify novel small ncRNA families. Here we present a **comprehensive** and **flexible** ncRNA analysis pipeline, **ncPRO-seq (Non-Coding RNA PRO**file from sRNA-seq), which is able to interrogate and perform detailed profiling analysis on small RNAs derived from annotated non-coding regions in **miRBase, Rfam** and **repeatMasker**, as well as regions defined by users. We perform both **gene-based** and **family-based** detailed analyses of small RNAs. The ncPRO-seq pipeline also has a module to identify regions significantly enriched with short reads that can not be classified as known ncRNA families [3], thus enabling the discovery of yet unknown ncRNA families. The ncPRO-seq pipeline supports input read sequences in fastq, fasta and color space format, as well as alignment results in BAM format, meaning that small RNA raw data from the 3 current major platforms (Roche-454, Illumina Solexa and Life technologies-SOLiD) could be analyzed with this pipeline. Finally, the ncPRO-seq pipeline can be used to analyze data based on genome from metazoan to plants. The current version proposes annotation files for **fifteen** different species.

### ncPRO-seq Workflow



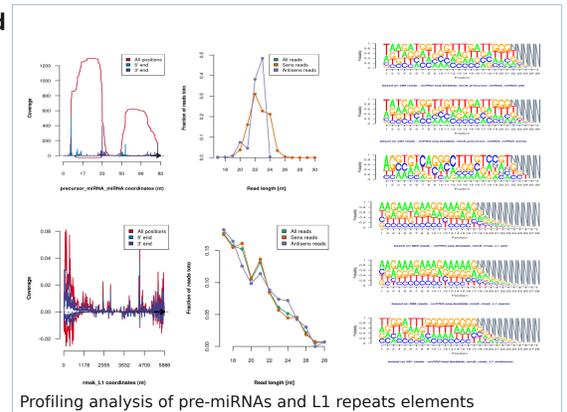
- ➔ Support multiple Solexa, SOLiD, 454 raw reads, and Bam files
- ➔ Reads grouping strategy (distinct vs abundant reads)
- ➔ Quality control of raw and aligned reads
- ➔ Reads mapping using the Bowtie software [4]
- ➔ More than 15 annotated organisms from mammals/metazoan to plants
- ➔ Flexible annotation and analysis of ncRNA families from Rfam, UCSC tRNA and miRBase
- ➔ Annotation and analysis of repeats classes from RepeatMasker
- ➔ Support user defined annotation files (gff3)
- ➔ Detect regions significantly enriched with reads
- ➔ Settings of UCSC Genome Browser tracks for visualization
- ➔ Stand-alone/command line pipeline and user-friendly interface

### Running ncPRO-seq



### Quality Control and mapping

### Family-based Analysis



### ncPRO-seq Interactive Analysis Report

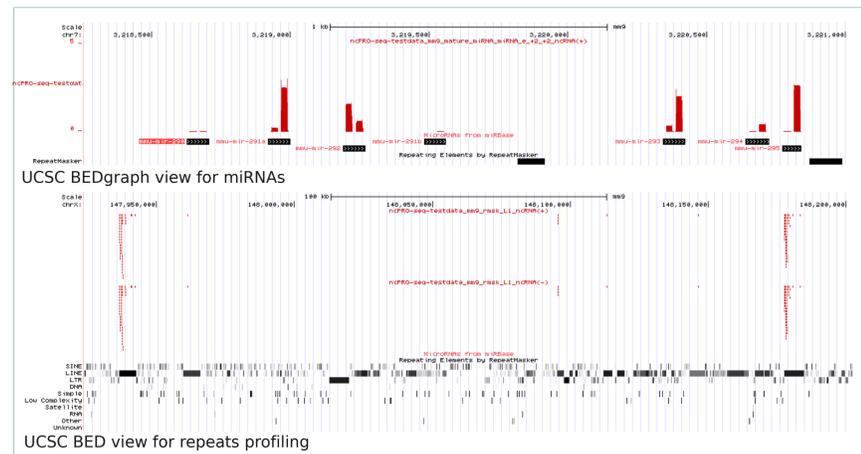
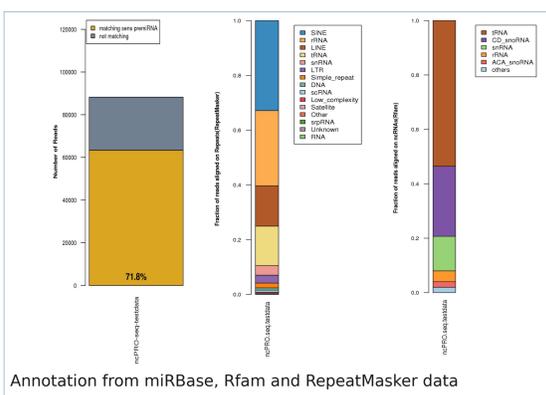


### Export and Visualization

Table view of annotation counts

Gene	RepeatMasker	Rfam	miRBase	ncPRO-seq
miR-101-2	1701	223333333		
miR-101-1	1000	1000		
miR-101-3	1000	1000		
miR-101-4	1000	1000		
miR-101-5	1000	1000		
miR-101-6	1000	1000		
miR-101-7	1000	1000		
miR-101-8	1000	1000		
miR-101-9	1000	1000		
miR-101-10	1000	1000		
miR-101-11	1000	1000		
miR-101-12	1000	1000		
miR-101-13	1000	1000		
miR-101-14	1000	1000		
miR-101-15	1000	1000		
miR-101-16	1000	1000		
miR-101-17	1000	1000		
miR-101-18	1000	1000		
miR-101-19	1000	1000		
miR-101-20	1000	1000		
miR-101-21	1000	1000		
miR-101-22	1000	1000		
miR-101-23	1000	1000		
miR-101-24	1000	1000		
miR-101-25	1000	1000		
miR-101-26	1000	1000		
miR-101-27	1000	1000		
miR-101-28	1000	1000		
miR-101-29	1000	1000		
miR-101-30	1000	1000		

### Annotation of ncRNAs and repeats classes



The ncPRO-seq pipeline is a stand-alone pipeline, which can be easily installed in a local computer or cluster. We offer two ways to launch the pipeline, through either a command line or a user-friendly web interface. The ncPRO-seq pipeline allows users to specify different options at each analysis stage, from raw reads processing to ways to generate results, all of which can be done by either selecting parameters in the web page or manually editing a configuration file. The results are available through an HTML report. Users can directly view figures and tables in the result web page. Track files are generated for visualization in genome browsers.

We deploy the ncPRO-seq pipeline in <http://ncproseq.sourceforge.net>, where users can find detailed information, such as basic descriptions, manuals, test dataset and example results. An online version for small dataset is available at <http://ncproseq.sourceforge.net/online.html>

[1] M. Ghildiyal, P.D. Zamore. Small silencing RNAs: an expanding universe. Nat Rev Genet., 10(2):94-108, 2009.  
 [2] P. Brodersen, O. Voinnet. The diversity of RNA silencing pathways in plants. Trends Genet., 22(5):268-80, 2006  
 [3] J. Toedling, C. Ciaudo, O. Voinnet, E. Heard, and E. Barillot. Girafe-an R/Bioconductor package for functional exploration of aligned next-generation sequencing reads. Bioinformatics, 26, 2902-2903, 2010.  
 [4] B. Langmead, C. Trapnell, et al. (2009). Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. Genome Biol., 10, R25.