

# Pipelines

## 1 Pipeline for genome decontamination (DeCon)

### 1.1 Introduction

Pipeline DeCon is designed to retrieve genomic sequences of target **phylum** from metagenomic assembly of paired next generation sequencing (NGS) reads. First, contigs below 400 base pair (bp) are removed and NGS reads are mapped to assembly by minimap2 (Li 2018), generating BAM file. Second, SprayNPray (Garber et al. 2022) is used to compute coverage, GC content and coding density of each contigs. Third, all contigs are searched against non-redundant (nr) database by DIAMOND (Buchfink et al. 2015) and assigned to phyla by MEGAN (Huson et al. 2007). Forth, a decision tree classifier is trained, taking coverage, GC content and coding density as training features and phylum assignment as target value. This classifier is used to compute phylum assignment of contigs that DIAMOND and MEGAN failed to compute assignments. Fifth, contigs assigned to the target phylum are retrieved. QUASt (Gurevich et al. 2013) and BUSCO (Simão et al. 2015) are used to evaluate retrieved genome. Distributions of contig coverage and GC content of retrieved genome are plotted.

### 1.2 Dependencies

#### Softwares

R

Python

minimap2

SAMtools

SprayNPray

DIAMOND

MEGAN (blast2rma rma2info scripts)

seqkit

QUAST

25 BUSCO

26

## 27 **Databases**

28 DIAMOND database (nr)

29 MEGAN database

30 BUSCO database

31

## 32 **Python modules**

33 numpy

34 pandas

35 scikit-learn

36

## 37 **R packages**

38 reticulate

39 stringr

40 ggplot2

41 ggExtra

42

## 43 **1.3 Usage**

44 Modify configuration file (templated as DeCon.conf.R), and run

45 Rscript path/DeCon\_pipeline.R path/DeCon\_main.R path/DeCon\_main.py path/DeCon.conf.R

# 46 **2 Pipeline for calling protein-coding genes from genome (Prot-** 47 **GeneCall)**

## 48 **2.1 Introduction**

49 Pipeline ProtGeneCall is designed to call protein-coding genes from genome, combining protein-genome  
50 alignments, transcriptome-genome alignments and *ab initio* gene predictions. First, repeat elements are  
51 identified by RepeatModeler (Smit et al. 2015b) and masked by RepeatMasker (Smit et al. 2015a). Masked  
52 genome is used for downstream analysis. Second, proteins of closely related species are mapped to the  
53 masked genome by miniprot (Li 2023). Third, paired RNA-sequencing (RNA-seq) reads are mapped to  
54 masked genome by Hisat2 (Kim et al. 2019). Forth, transcriptome-genome alignments are computed by

StringTie (Pertea et al. 2015). Fifth, gene structures are predicted from transcriptome-genome alignments by TransDecoder (Haas et al. 2016), combining searching against UniRef and PfamA databases. Sixth, AUGUSTUS (Stanke et al. 2003) is trained with gene structures from TransDecoder to compute gene predictions. Seventh, BRAKER (Hoff et al. 2019) is trained with RNA-seq mapping to call genes. Eighth, GALBA (Hoff et al. 2019) is trained with proteins of closely related species. Ninth, protein-genome alignments from miniprot, transcript-genome alignments from Hisat2-StringTie, and *ab initio* gene predictions from TransDecoder, AUGUSTUS, BRAKER and GALBA are integrated into consensus gene structures by EvidenceModeler (Haas et al. 2008). Tenth, genes supported by only one *ab initio* predictor and lack protein/RNA-seq evidence are removed. Eleventh, PASA (Haas et al. 2008) is run twice to update filtered gene structures from EvidenceModeler. Twelfth, genes with in-frame stop codons are removed and the predicted peptide set is evaluated by BUSCO (Simão et al. 2015).

## 2.2 Dependencies

### Softwares

R  
Python  
RepeatModeler  
RepeatMasker  
miniprot  
Hisat2  
SAMtools  
StringTie  
TransDecoder  
HMMER  
DIAMOND  
AGAT  
AUGUSTUS  
BLAST+  
GALBA  
BRAKER  
EvidenceModeler  
BUSCO  
gffread

87 seqkit

88 MAKER

89

## 90 **Databases**

91 DIAMOND database (UniRef)

92 Pfam-A

93

## 94 **External scripts**

95 cufflinks\_gtf\_to\_alignment\_gff3.pl from EvidenceModeler

96 augustus\_GFF3\_to\_EVM\_GFF3.pl from EvidenceModeler

97 gth2gtf.pl from AUGUSTUS

98 computeFlankingRegion.pl from AUGUSTUS

99 gff2gbSmallDNA.pl from AUGUSTUS

100 gtf2aa.pl from AUGUSTUS

101 simplifyFastaHeaders.pl from AUGUSTUS

102 aa2nonred.pl from AUGUSTUS

103 filterGenesIn.pl from AUGUSTUS

104 autoAug.pl from AUGUSTUS

105 evm\_evidence.py in this GitHub

106

## 107 **R packages**

108 stringr

109 parallel

110

## 111 **2.3 Usage**

112 Modify configuration file (templated as ProtGeneCall\_conf.R), and run

113 Rscript path/ProtGeneCall\_pipeline.R path/ProtGeneCall\_main.R path/ProtGeneCall\_conf.R

## 114 **3 Pipeline for calling repeat elements from genome (RepCall)**

### 115 **3.1 Introduction**

116 Pipeline RepCall is designed to call repeat elements genes from genome. First, miniature inverted-repeat  
117 transposable elements (MITE) are called by MITE-Hunter (Han et al. 2010). Second, long terminal re-  
118 peats (LTRs) are identified by incorporating LTR\_FINDER\_parallel (Ou et al. 2019), LTRharvest (Elling-  
119 haus et al. 2008) and LTR\_retriever (Ou et al. 2018). Third, identified MITEs and LTRs are masked by  
120 RepeatMasker (Smit et al. 2015a). Forth, RepeatModeler (Smit et al. 2015b) is used to further identify  
121 repeats in the masked genome. Fifth, the locations of MITEs, LTRs and repeats from RepeatModeler are  
122 identified by RepeatMasker and all repeats are incorporated into a consensus library.

### 123 **3.2 Dependencies**

#### 124 **Softwares**

125 R  
126 seqkit  
127 MITE-Hunter  
128 LTR\_FINDER\_parallel  
129 LTRharvest  
130 LTR\_retriever  
131 RepeatMasker  
132 RepeatModeler

### 134 **3.3 Usage**

135 Modify configuration file (templated as RepCall\_conf.R), and run  
136 Rscript path/RepCall\_pipeline.R path/RepCall\_main.R path/RepCall\_conf.R

## 137 **4 Pipeline for calling non-coding RNA (ncRNAcall)**

### 138 **4.1 Introduction**

139 Pipeline ncRNAcall is designed to call non-coding RNA (ncRNA) from genome. First, transfer RNA  
140 (tRNA) is identified by tRNAscan-SE (Lowe et al. 1997). Second, microRNA is called by miRNature  
141 (Velandia-Huerto et al. 2021). Third, target genes of microRNA are identified by searching microRNA

142 against annotated three prime untranslated regions (3'UTR) by miRanda (Enright et al. 2003). Forth, In-  
143 fernal (Nawrocki et al. 2013) searches against Rfam (Kalvari et al. 2021) database to call other non-coding  
144 RNA, *e.g.* ribosomal RNA (rRNA) and small nuclear RNA (snRNA). Fifth, all results are incorporated  
145 together.

## 146 4.2 Dependencies

### 147 Softwares

148 R  
149 tRNAscan-SE  
150 biocode  
151 miRNature  
152 miRanda  
153 bedtools  
154 seqkit  
155 Infernal

### 156 Databases

157 miRNature database  
158 Rfam database

### 160 R packages parallel

161 stringr

## 164 4.3 Usage

165 Modify configuration file (templated as ncRNAcall.conf.R), and run  
166 Rscript path/ncRNAcall\_pipeline.R path/ncRNAcall\_main.R path/ncRNAcall.conf.R

## 167 References

168 Buchfink, Benjamin et al. (2015). “Fast and sensitive protein alignment using DIAMOND”. In: *Nature*  
169 *methods* 12.1, pp. 59–60.  
170 Ellinghaus, David et al. (2008). “LTRharvest, an efficient and flexible software for de novo detection of  
171 LTR retrotransposons”. In: *BMC bioinformatics* 9, pp. 1–14.

172 Enright, Anton et al. (2003). “MicroRNA targets in *Drosophila*”. In: *Genome biology* 4, pp. 1–27.

173 Garber, Arkadiy I et al. (2022). “SprayNPray: user-friendly taxonomic profiling of genome and  
174 metagenome contigs”. In: *BMC genomics* 23.1, p. 202.

175 Gurevich, Alexey et al. (2013). “QUAST: quality assessment tool for genome assemblies”. In: *Bioinfor-*  
176 *matics* 29.8, pp. 1072–1075.

177 Haas, B et al. (2016). “TransDecoder (find coding regions within transcripts)”. In: *Google Scholar*.

178 Haas, Brian J et al. (2008). “Automated eukaryotic gene structure annotation using EVIDENCEModeler  
179 and the Program to Assemble Spliced Alignments”. In: *Genome biology* 9, pp. 1–22.

180 Han, Yujun et al. (2010). “MITE-Hunter: a program for discovering miniature inverted-repeat transposable  
181 elements from genomic sequences”. In: *Nucleic acids research* 38.22, e199–e199.

182 Hoff, Katharina J et al. (2019). “Whole-genome annotation with BRAKER”. In: *Gene prediction: methods*  
183 *and protocols*, pp. 65–95.

184 Huson, Daniel H et al. (2007). “MEGAN analysis of metagenomic data”. In: *Genome research* 17.3,  
185 pp. 377–386.

186 Kalvari, Ioanna et al. (2021). “Rfam 14: expanded coverage of metagenomic, viral and microRNA families”.  
187 In: *Nucleic Acids Research* 49.D1, pp. D192–D200.

188 Kim, Daehwan et al. (2019). “Graph-based genome alignment and genotyping with HISAT2 and HISAT-  
189 genotype”. In: *Nature biotechnology* 37.8, pp. 907–915.

190 Li, Heng (2018). “Minimap2: pairwise alignment for nucleotide sequences”. In: *Bioinformatics* 34.18,  
191 pp. 3094–3100.

192 — (2023). “Protein-to-genome alignment with minimot”. In: *Bioinformatics* 39.1, btad014.

193 Lowe, Todd M et al. (1997). “tRNAscan-SE: a program for improved detection of transfer RNA genes in  
194 genomic sequence”. In: *Nucleic acids research* 25.5, pp. 955–964.

195 Nawrocki, Eric P et al. (2013). “Infernal 1.1: 100-fold faster RNA homology searches”. In: *Bioinformatics*  
196 29.22, pp. 2933–2935.

197 Ou, Shujun et al. (2018). “LTR\_retriever: a highly accurate and sensitive program for identification of  
198 long terminal repeat retrotransposons”. In: *Plant physiology* 176.2, pp. 1410–1422.

199 — (2019). “LTR\_FINDER\_parallel: parallelization of LTR\_FINDER enabling rapid identification of long  
200 terminal repeat retrotransposons”. In: *Mobile DNA* 10.1, pp. 1–3.

201 Pertea, Mihaela et al. (2015). “StringTie enables improved reconstruction of a transcriptome from RNA-  
202 seq reads”. In: *Nature biotechnology* 33.3, pp. 290–295.

203 Simão, Felipe A et al. (2015). “BUSCO: assessing genome assembly and annotation completeness with  
204 single-copy orthologs”. In: *Bioinformatics* 31.19, pp. 3210–3212.

205 Smit, AFA et al. (2015a). *RepeatMasker Open-4.0. 2013–2015*.

206 Smit, AFA et al. (2015b). “RepeatModeler Open-1.0. 2008–2015”. In: *Seattle, USA: Institute for Systems*  
 207 *Biology. Available from: [httpwww.repeatmasker.org](http://www.repeatmasker.org), Last Accessed May 1*, p. 2018.

208 Stanke, Mario et al. (2003). “Gene prediction with a hidden Markov model and a new intron submodel”.  
 209 In: *Bioinformatics* 19.suppl\_2, pp. ii215–ii225.

210 Velandia-Huerto, Cristian A et al. (2021). “miRNAture—Computational Detection of microRNA Candi-  
 211 dates”. In: *Genes* 12.3, p. 348.