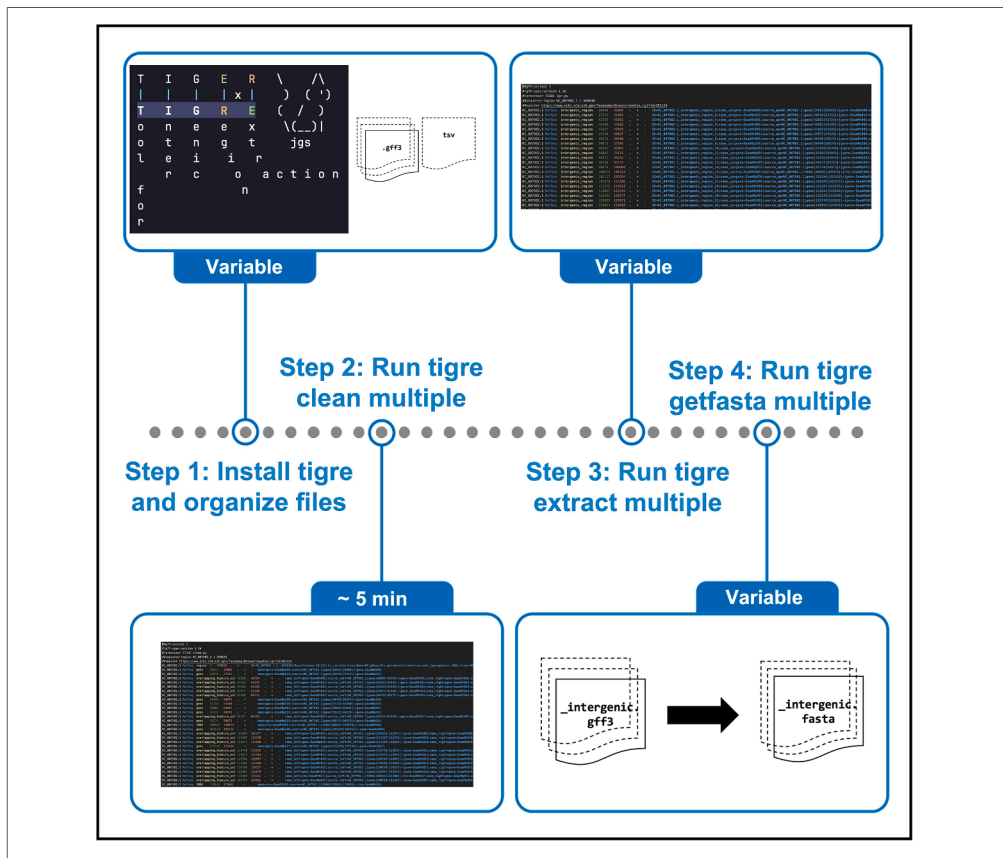


## Protocol

# Protocol for extracting intergenic regions from annotated genomes using TIGRE



Here, we present a protocol for extracting intergenic regions from annotated genomes using TIGRE (Tool for Intergenic Region Extraction). This protocol describes steps for installing and running TIGRE locally. Via TIGRE, the user can prepare the input data ("clean" command), retrieve the intergenic sequences ("extract" command), obtain the respective intergenic nucleotide sequences ("getfasta" command), and merge different GFF3 files ("combine" command). TIGRE capabilities are demonstrated on the retrieval of intergenic sequences from 1,207 land plant mitochondrial genomes.

Publisher's note: Undertaking any experimental protocol requires adherence to local institutional guidelines for laboratory safety and ethics.

Breno Dupin,  
Matheus Sanita  
Lima, Alexandre  
Rossi Paschoal,  
David Roy Smith

breno.dupin@gmail.com  
(B.D.)  
msanital@uwo.ca (M.S.L.)

**Highlights**  
Protocol for  
extracting intergenic  
regions from  
annotated genomes

Instructions on how to  
interact with TIGRE's  
CLI (command-line  
interface)

Guidance on how to  
implement the  
advanced use of  
TIGRE for the  
extraction of introns

Dupin et al., STAR Protocols 6,  
104270  
December 19, 2025 © 2025  
The Author(s). Published by  
Elsevier Inc.  
<https://doi.org/10.1016/j.xpro.2025.104270>



## Protocol

## Protocol for extracting intergenic regions from annotated genomes using TIGRE

Breno Dupin,<sup>1,4,5,\*</sup> Matheus Sanita Lima,<sup>2,4,6,\*</sup> Alexandre Rossi Paschoal,<sup>1,3</sup> and David Roy Smith<sup>2</sup><sup>1</sup>Department of Computer Science, Federal University of Technology-Parana, Cornelio Procopio, Paraná 86300-000, Brazil<sup>2</sup>Department of Biology, Western University, London, ON N6A 5B7, Canada<sup>3</sup>Rosalind Franklin Institute, Harwell Science Campus, Didcot OX11 0QX, UK<sup>4</sup>These authors contributed equally<sup>5</sup>Technical contact<sup>6</sup>Lead contact\*Correspondence: [breno.dupin@gmail.com](mailto:breno.dupin@gmail.com) (B.D.), [msanital@uwo.ca](mailto:msanital@uwo.ca) (M.S.L.)  
<https://doi.org/10.1016/j.xpro.2025.104270>

## SUMMARY

Here, we present a protocol for extracting intergenic regions from annotated genomes using TIGRE (Tool for Intergenic Region Extraction). This protocol describes steps for installing and running TIGRE locally. Via TIGRE, the user can prepare the input data ("clean" command), retrieve the intergenic sequences ("extract" command), obtain the respective intergenic nucleotide sequences ("getfasta" command), and merge different GFF3 files ("combine" command). TIGRE capabilities are demonstrated on the retrieval of intergenic sequences from 1,207 land plant mitochondrial genomes.

## BEFORE YOU BEGIN

In molecular evolution, few topics are more controversial than noncoding DNA. The two camps (functional vs junk DNA) have been at loggerheads for decades.<sup>1</sup> Between the catchy article titles<sup>2</sup> and hyperbolic media coverage,<sup>3</sup> three facts are settled: i) population genetics can easily explain the accumulation of noncoding DNA via nonadaptive mechanisms<sup>4</sup>; ii) pervasive transcription, although largely spurious,<sup>5</sup> is found in multiple systems<sup>6,7</sup>; and iii) noncoding RNAs have essential regulatory roles across the Tree of Life.<sup>7-9</sup>

With controversy comes scrutiny, and dozens of research groups have peered through noncoding nucleotides of prokaryotes and eukaryotes alike. Of particular interest is the evolution of intergenic regions (also named intergenic spacers or intergenic sequences), which make up most of the non-coding portion of any genome.<sup>10</sup> Most studies of intergenic regions gloss over the methods used to extract such sequences, but bioinformatics tools and protocols dedicated (exclusively or not) to the extraction of intergenic regions have been published.<sup>9,11-14</sup>

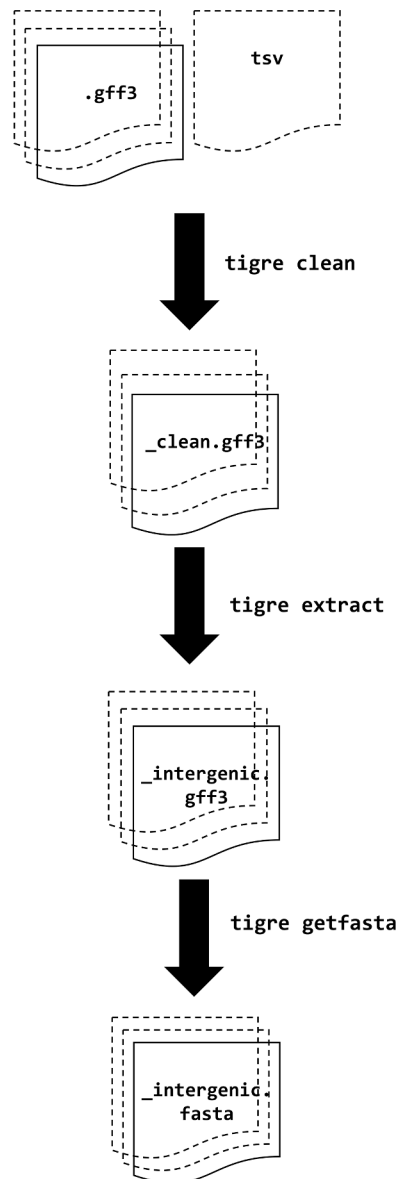
## Innovation

With a user-centered focus, TIGRE has been designed with scalability, interoperability, and customizability in mind. Contrary to traditional intergenic extraction strategies (such as bedtools complement), TIGRE maintains the identifying information of the delimiting genes and handles circular genomes. For even greater consistency and interoperability, the user can implement the Gene Dictionary Tool (GDT) and protocol<sup>15</sup> to standardize gene names across the study sample. This protocol facilitates the extraction and analysis of intergenic regions from any type of annotated genome.

## Overview

TIGRE - Tool for InterGenic Region Extraction – can extract intergenic regions from any annotated genome in a scalable manner. The tool has three main functions: clean, extract, and





**Figure 1. Overview of the TIGRE protocol**

TIGRE's pipeline is composed of three main commands - clean, extract, and getfasta. All commands can operate on a single or multiple GFF3 files. When processing multiple genomes, the genome annotations must be indexed in a TSV file. `tigre clean` does the "heavy lifting", as it standardizes the input GFF3 files for the extraction of intergenic regions. `tigre extract` generates GFF3 files containing the extracted intergenic regions, including regions spanning the genome boundaries of circular genomes. `tigre getfasta` retrieves the corresponding nucleotide sequences of the intergenic regions. `tigre combine` is an auxiliary command, so it is not depicted in the main pipeline. Details on the implementation of `tigre combine` can be found under Advanced Usage in <https://github.com/brenodupin/tigre>.

`getfasta` (Figure 1 and [step-by-step method details](#)). The main inputs are a TSV file indexing the genome(s) of interest and their corresponding GFF3 files. Should the user want to extract the nucleotide sequences of the intergenic regions, then the respective FASTA files of each genome are also necessary. TIGRE has a fourth auxiliary command called `combine`. This command is used only in the Advanced Usage of TIGRE for the extraction of introns (see [expected outcomes](#)).

**Table 1. Performance metrics for tigre clean, extract, and getfasta in two machines: Ubuntu 22.04.5 LTS and macOS Sonoma 14.6.1**

Command	Operating system	Elapsed time (s)	Max memory use (MB)
clean	Ubuntu 22.04.5 LTS	7.88	144.8
	macOS Sonoma 14.6.1	27.05	156.2
extract	Ubuntu 22.04.5 LTS	1.46	104.6
	macOS Sonoma 14.6.1	6.93	93
getfasta	Ubuntu 22.04.5 LTS	1.38	106.3
	macOS Sonoma 14.6.1	5.67	125.3

All commands share a unified logging system and two modes: simple vs multiple. The log flags (-v, -log, -no-log-file, -quiet) give the user autonomy to determine the level of detail to be logged in any TIGRE command. Unless -no-log-file is passed, a log will always be saved, at least according to the default parameters.

For added safety, TIGRE has checkpoints that ensure whether the input (and output) files already exist. Unless the flag -overwrite (which is shared across all TIGRE commands) is specified, output files will not be overwritten on already existing files with the same name. Just as defensively, TIGRE commands will create empty output files when the input files are empty. For instance, tigre clean might create a clean GFF containing overlapping feature sets juxtaposed with other genome features that leave no room for intergenic regions. In that case, tigre extract will generate an empty GFF (containing only the respective file header), and tigre getfasta will generate an empty file (without any header). These empty files function as a “cushion” preventing the program from crashing when the checkpoints (listed above) are called.

TIGRE’s potential for scalability is inherently designed in its built-in multiprocessing capabilities. The modes single and multiple allow the user to choose between running a quick test in one genome (single) or a large-scale analysis in thousands of genomes (multiple). TIGRE’s multiprocessing framework also supports the integration of GDT,<sup>15</sup> regardless of the size of the gene dictionary files (.gdict) and/or the number of analyzed genomes (Table 1).

The three commands are intertwined via a group of shared flags that allows the user to create customized pipelines (Figure 2). Depending on the genome features specified via the -query-string flag (in tigre clean), the user can extract not only intergenic regions but any region(s) that falls between any set(s) of genome feature(s). For example, if the user specifies -query-string as being “type in (‘tRNA’, ‘region)’”, tigre clean will filter out all other genome features (such as genes, rRNAs, etc.) and provide tigre extract with a clean GFF file that has only tRNAs as the delimiting elements of the soon-to-be-extracted “inter-tRNA” regions. The several -gff-<> and -fasta-<> flags, combined with -feature-type, give users the flexibility to customize their pipeline according to their interest.

Optional flags such as -keep-orfs (for tigre clean), -add-region (for tigre extract), and -bedtools-compatible (for tigre getfasta) were designed to meet the needs and challenges usually encountered when analyzing organelle intergenic regions across disparate species. Detailed documentation of all three commands and their multiple flags, can be found at <https://github.com/brenodupin/tigre>.

### Download the software and prerequisites

TIGRE is a pure Python (>=3.12) CLI tool developed in Unix environment. TIGRE was tested extensively on Ubuntu and validated for use on MacOS. The tool can be downloaded and installed directly from PyPI (<https://pypi.org/project/tigre/>). The only mandatory dependency is pandas. Biopython and GDT are the two other optional dependencies: tigre getfasta relies on Biopython to extract FASTA sequences of the intergenic regions; and tigre clean implements GDT, if -gdict is used, to standardize gene names when formatting its output (<AN>\_clean.gff3).

```

) tigre --help
usage: tigre [-h] [-v] [--log PATH] [--no-log-file] [--quiet] [--version]
           {clean,extract,getfasta,combine} ...

  T I G E R \ \
  | | | | x | ) ( '
  T I G R E ( / )
  o n e e x \(__|
  o t n g t jgs
  l e i i r
  r c o a c t i o n
  f
  n
  o
  r

positional arguments:
  {clean,extract,getfasta,combine}
  clean                Clean, standardize, and prepare GFF3 file(s) for
                        processing.
  extract              Extract intergenic regions from processed GFF3
                        file(s).
  getfasta             Generate FASTA sequences from intergenic regions.
  combine              Combine two GFF3 files into one.

options:
  -h, --help          show this help message and exit
  --version           Show the version of the tigre package.

log options:
  -v, --verbose       Increase verbosity level. Use multiple times for more
                        verbose output: -v (INFO console, TRACE file), -vv
                        (DEBUG console, TRACE file), -vvv (TRACE console,
                        TRACE file). Default: INFO console + DEBUG file.
  --log PATH          Path to the log file. If not provided, a default log
                        file will be created.
  --no-log-file       Disable file logging.
  --quiet             Suppress console output.

) tigre clean --help
usage: tigre clean [-h] [-v] [--log PATH] [--no-log-file] [--quiet] [--gdict PATH] [--query-string STR] [--keep-orfs] [--overwrite] [--extended-filtering] mode ...

positional arguments:
  mode
  single                Process a single GFF3 file
  multiple             Process multiple GFF3 files from TSV

options:
  -h, --help          show this help message and exit

log options:
  -v, --verbose       Increase verbosity level. Use multiple times for more verbose output: -v (INFO console, TRACE file), -vv
                        (TRACE console, TRACE file), -vvv (DEBUG console, TRACE file), -vvv
                        (DEBUG console, TRACE file). Default: INFO console + DEBUG file.
  --log PATH          Path to the log file. If not provided, a default log file will be created.
  --no-log-file       Disable file logging.
  --quiet             Suppress console output.

clean options:
  --gdict PATH        Path to a Gene Dictionary Table (GDT) file to standardize gene names. Optional.
  --query-string STR pandas query string to filter features in GFF3 file. Default: 'type in ('gene', 'trna', 'rRNA', 'region')'.
  --keep-orfs         Keep ORF sequences in the output. Default: False.
  --overwrite         Overwrite existing output files. Default: False.
  --extended-filtering Enable extended filtering to more aggressively remove ORFs and their related features. This may remove more features than intended in some
                        cases. Default: False.

) tigre extract --help
usage: tigre extract [-h] [-v] [--log PATH] [--no-log-file] [--quiet] [--add-region] [--overwrite] [--feature-type STR] mode ...

extract options:
  --add-region        Add region line to the output GFF3 file.
  --overwrite         Overwrite existing output files.
  --feature-type STR Feature type name to use for intervening regions in the output GFF3 file. Features spanning the genome boundary in circular genomes will be
                        appended with '_merged'. Default: 'intergenic_region'.

) tigre getfasta --help
usage: tigre getfasta [-h] [-v] [--log PATH] [--no-log-file] [--quiet] [--bedtools-compatible] [--overwrite] mode ...

getfasta options:
  --bedtools-compatible Adjust FASTA headers to use 0-based indexing for bedtools compatibility. Affects headers only, not sequences. Default: False.
  --overwrite           Overwrite existing output files. Default: False.

) tigre combine --help
usage: tigre combine [-h] [-v] [--log PATH] [--no-log-file] [--quiet] [--overwrite] mode ...

combine options:
  --overwrite         Overwrite existing output files. Default: False.

```

Figure 2. TIGRE's command-line interface (CLI)

TIGRE's command-line interface is extensively documented to help users make sense of the pipeline. Further implementation details of each command can be found at <https://github.com/brenodupin/tigre>. Detailed

### Figure 2. Continued

explanation of how TIGRE handles circular genomes is presented in <https://github.com/brenodupin/tigre/blob/master/PROCESSING.md>. The logo contains the ASCII art from Joan G. Stark ([https://oldcompcz.github.io/jgs/joan\\_stark/index-2.html](https://oldcompcz.github.io/jgs/joan_stark/index-2.html)).

**Note:** TIGRE's built-in optimization strategies allow for a wide range of hardware specifications. Any modern computer (desktop or laptop) with at least 8 GB of RAM should be able to use TIGRE for analyzing thousands of organelle genomes. For reference, TIGRE's performance was tested using a dataset with 1,207 land plant mitochondrial genomes in two machines: Ubuntu 22.04.5 LTS laboratory server (dual Intel Xeon E5-2640 v4 @ 2.40 GHz; 20 cores / 40 threads); 384 GB RAM, and macOS Sonoma 14.6.1 (Intel Core i9-8950HK @ 2.90 GHz; 6 cores / 12 threads; 32 GB RAM) (Table 1). All tests used the maximal available number of threads.

## KEY RESOURCES TABLE

REAGENT or RESOURCE	SOURCE	IDENTIFIER
<b>Deposited data</b>		
Updated gene dictionary for land plant mitochondrial genomes (.gdt) files	Dupin et al. <sup>15</sup>	<a href="https://zenodo.org/records/17177240">https://zenodo.org/records/17177240</a>
Original GFF3 and FASTA files of 1.2K land plant mitochondrial genomes	NCBI Nucleotide	<a href="https://www.ncbi.nlm.nih.gov/nucleotide/">https://www.ncbi.nlm.nih.gov/nucleotide/</a> and <a href="https://zenodo.org/records/17475986">https://zenodo.org/records/17475986</a>
Output GFF3 and FASTA files of 1.2K land plant mitochondrial genomes	This manuscript	<a href="https://zenodo.org/records/17475986">https://zenodo.org/records/17475986</a>
<b>Software and algorithms</b>		
gdt library	Dupin et al. <sup>15</sup>	<a href="https://github.com/brenodupin/gdt">https://github.com/brenodupin/gdt</a>
Python 3.12	Python Software Foundation	<a href="https://www.python.org">https://www.python.org</a>
Biopython 1.85	Cock et al. <sup>16</sup>	<a href="https://biopython.org/">https://biopython.org/</a>
pandas 2.2.3	The Pandas Development Team	<a href="https://pandas.pydata.org/">https://pandas.pydata.org/</a>

## MATERIALS AND EQUIPMENT

As a pure Python CLI tool, TIGRE has a user-friendly interface, permitting users to implement flexibly the different commands (clean, extract, getfasta, and combine). TIGRE's functionalities are organized in their own Python files: cli.py is the actual CLI, so it orchestrates all tigre commands and arguments passed by the user; gff3\_utils.py is derived from the gdt library and possess several functionalities to treat GFF files; log\_setup.py organizes all the detailed logging system of TIGRE; clean.py, and clean\_gdt.py make up the tigre clean command in its various flavours (i.e., single vs multiple and with or without GDT); biopython\_wrapper.py is the script for the tigre getfasta command; igr.py is the cornerstone of TIGRE, as this file encapsulates the intergenic region extraction rationale and mechanics; combine.py contains the source code for the auxiliary command that merges GFF3 files. The annotated genomes must be inputted as GFF3 files. The mandatory and optional dependencies are listed in the [key resources table](#). The source code for TIGRE can be found at <https://github.com/brenodupin/tigre>.

## STEP-BY-STEP METHOD DETAILS

### Preparing the environment and files

⌚ **Timing:** Minutes to hours (depending on dataset size)

This section sets up the environment and organizes the input data for the extraction of intergenic sequences.

1. Install TIGRE, and possibly biopython and gdt, using pip:

```
>pip install "tigre[all]"
```

**Note:** "tigre[all]" will automatically install the two optional dependencies gdt and biopython, which are needed to analyze the example dataset. These extras can be installed separately with "tigre[gdt]" and "tigre[bio]". You can also install tigre without any optional dependency via ">pip install tigre".

**Note:** There is an example dataset in <https://github.com/brenodupin/tigre/tree/master/examples>. If you would like to test this protocol using the example dataset (along with the commands below), you should download the plants\_mit.tar.zst (compressed) file and decompress it with >tar -xf plants\_mit.tar.zst. Then you should enter the folder with >cd plants\_mit and run TIGRE from it. If you follow these steps, you can skip to step 4 (See [troubleshooting](#) for details on how to use the example dataset for benchmarking).

2. Create a TSV file listing the accession numbers of the genomes you want to analyze.
3. Create a separate folder for each genome listed in the TSV, in the working directory containing your TSV file.

**Note:** The example dataset contains only annotations from NCBI, so the TSV file lists NCBI Accession Numbers (as in NC\_XXXX.X). Each folder must be named exactly as the corresponding genome annotation listed in the TSV - e.g., folder as "NC\_007982.1", and inside it the GFF3 file as "NC\_007982.1.gff3".

⚠ **CRITICAL:** The column in the TSV file must be named. The (default) name used in the example dataset is "AN" (for Accession Number), but you can name this column in any other way. The correct column name must be passed to the `-an-column` flag, whenever this flag is needed. Similarly, the GFF files can have a file extension different than `.gff3`; the correct extension must be passed to `-gff-in-ext`. If you wish to maintain consistency across your dataset, the flag `-gff-out-ext` can be used in certain commands to produce output files with the same extension. Note that TIGRE does not convert any file types - i.e., the GFF files will always be GFF files regardless of the chosen extension.

### Running tigre clean

⌚ **Timing:** ~5 min (depending on dataset size)

This section standardizes the inputted GFF3 files in a way that allows the extraction of intergenic regions from the genomes of interest.

4. Run `tigre clean multiple` from your working direct:

```
>tigre clean multiple -v --log tigre_clean.log --gdict plants_mit.gdict --tsv plants_mit.tsv --overwrite
```

**Note:** In the example commands, the flags `-v`, `-log`, and `-gdict` are used to showcase TIGRE's functionalities. These flags are optional and can be used independent of each other. The `-v` flag will provide the user with details for later debugging, `-gdict` implements the gene

dictionary capabilities from GDT, and `-log` will point to the file where the log will be saved. As default the log is saved in the working directory as `tigre_<command>_<timestamp>.log`. There also exist `-quiet` and `-no-log-file`; the first one will not print any message onto stdout (standard output), whereas the second one will not create a log file.

**Note:** `tigre clean` will do the heavy lifting in TIGRE. This command will standardize the GFF files, so they can serve as input for the next commands. The standardization happens in five main ways: i) filtering of chosen genome features that have been determined by the user via the `-query-string` flag; ii) genome features are sorted (in ascending order) according to their start coordinates; iii) overlapping features are merged into `"overlapping_feature_set"` while retaining the identifying information of their delimiting components; iv) genome boundaries are treated in a way that allows `tigre extract` and `tigre getfasta` to obtain intergenic regions that might be wrapping around the start/end region of the genome (<https://github.com/brenodupin/tigre/blob/master/PROCESSING.md>); v) the names and positions of the delimiting genes are maintained in the attributes column for traceability. By default, gene names are retrieved from the `"ID="` field in the attributes column of the original GFF file. Should the user want to standardize gene names across genomes, they can implement GDT - Gene Dictionary Tool - in this step. TIGRE supports gene dictionary files (`.gdict`) from the `gdt` library, and it has been vastly optimized for large input datasets and gene dictionaries.

△ **CRITICAL:** As an output, `tigre clean` will generate "clean" GFF3 files that contain only the genome features specified in the `-query-string`. In the example command, this flag is not passed, which means that the default value (`"type in ('gene', 'tRNA', 'rRNA', 'region')"`) is used. You can change this value to include your genome features of interest. If you are not using the default of `-query-string`, make sure to maintain `'region'` in your new `-query-string` list - e.g., `-query-string "type in ('gene', 'region')"`. The genome features that have remained in the "clean" GFF will serve as signposts representing the genome regions that codify something (typically proteins). The subsequent steps of `tigre` will extract everything that falls outside (or in between) these sign posts (Figure 3).

### Running `tigre extract`

⌚ **Timing:** Minutes (depending on the size of the dataset)

This section retrieves the intergenic regions from the "clean" GFF3 files. The intergenic regions are not exactly "extracted" but are "annotated" as intergenic regions.

5. Run `tigre extract` to obtain the intergenic regions from the clean GFF3 files:

```
>tigre extract multiple -v --log tigre_extract.log --tsv plants_mit.tsv --overwrite
```

**Note:** The output of `tigre extract` is the `"<AN>_intergenic.gff3"` files. These files are the "negative" of the clean GFF3 files. In the attributes column, each intergenic region is annotated with the identifying information of its delimiting genes (Figure 4). All intergenic regions are annotated as having positive polarity and they receive the name "intergenic region" under the type column.

△ **CRITICAL:** By default, the `"region"` element is not added to `<AN>_intergenic.gff3`. If you want to add this element to your intergenic GFF3 files, you should pass the `-add-region` flag.

```

##gff-version 3
##gff-spec-version 1.26
##processor TIGRE_clean.py
##sequence-region NC_007982.1 1 589800
##species https://www.ncbi.nlm.nih.gov/Taxonomy/Browser/wwwtax.cgi?l=381124
NC_007982.1 RefSeq region 1 589800 + . ID=NC_007982.1.1.589838;dbxref=taxon:381124;is_circular=true;Name=HT;gbkey=Src;genome=mitochondrion;mol_type=genomic DNA;strain=H
NC_007982.1 RefSeq gene 19813 28889 + . name=gene-ZeamP088;source=NC_007982.1|gene|19813|28889|+|gene-ZeamP088|
NC_007982.1 RefSeq gene 25918 27232 + . name=gene-ZeamP012;source=NC_007982.1|gene|25918|27232|+|gene-ZeamP012|
NC_007982.1 RefSeq overlapping_feature_set 42663 42734 + . name_left=gene-ZeamM081;source_left=NC_007982.1|gene|42663|42734|+|gene-ZeamM081;name_right=gene-ZeamM081;source_right=NC_007982.1|gene|42663|42734|+|gene-ZeamM081;name_right=gene-ZeamM081|
NC_007982.1 RefSeq overlapping_feature_set 43362 43435 + . name_left=gene-ZeamM082;source_left=NC_007982.1|gene|43362|43435|+|gene-ZeamM082;name_right=gene-ZeamM082;source_right=NC_007982.1|gene|43362|43435|+|gene-ZeamM082;name_right=gene-ZeamM082|
NC_007982.1 RefSeq overlapping_feature_set 44163 44236 + . name_left=gene-ZeamM083;source_left=NC_007982.1|gene|44163|44236|+|gene-ZeamM083;name_right=gene-ZeamM083;source_right=NC_007982.1|gene|44163|44236|+|gene-ZeamM083;name_right=gene-ZeamM083|
NC_007982.1 RefSeq overlapping_feature_set 4707 47148 + . name_left=gene-ZeamM084;source_left=NC_007982.1|gene|4707|47148|+|gene-ZeamM084;name_right=gene-ZeamM084;source_right=NC_007982.1|gene|4707|47148|+|gene-ZeamM084;name_right=gene-ZeamM084|
NC_007982.1 RefSeq overlapping_feature_set 49198 49272 + . name_left=gene-ZeamM085;source_left=NC_007982.1|gene|49198|49272|+|gene-ZeamM085;name_right=gene-ZeamM085;source_right=NC_007982.1|gene|49198|49272|+|gene-ZeamM085;name_right=gene-ZeamM085|
NC_007982.1 RefSeq gene 58490 58874 + . name=gene-ZeamP186;source=NC_007982.1|gene|58490|58874|+|gene-ZeamP186|
NC_007982.1 RefSeq gene 57194 59340 + . name=gene-ZeamP019;source=NC_007982.1|gene|57194|59340|+|gene-ZeamP019|
NC_007982.1 RefSeq gene 63968 64881 + . name=gene-ZeamP021;source=NC_007982.1|gene|63968|64881|+|gene-ZeamP021|
NC_007982.1 RefSeq gene 76132 84976 + . name=gene-ZeamP032;source=NC_007982.1|gene|76132|84976|+|gene-ZeamP032|
NC_007982.1 RefSeq overlapping_feature_set 94257 94329 + . name_left=gene-ZeamM087;source_left=NC_007982.1|gene|94257|94329|+|gene-ZeamM087;name_right=gene-ZeamM087;source_right=NC_007982.1|gene|94257|94329|+|gene-ZeamM087;name_right=gene-ZeamM087|
NC_007982.1 RefSeq gene 95774 98873 + . name=gene-ZeamP033;source=NC_007982.1|gene|95774|98873|+|gene-ZeamP033|
NC_007982.1 RefSeq rRNA 180893 180978 + . name=rna-ZeamM088;source=NC_007982.1|rRNA|180893|180978|+|rna-ZeamM088|
NC_007982.1 RefSeq gene 181544 182116 + . name=gene-ZeamP036;source=NC_007982.1|gene|181544|182116|+|gene-ZeamP036|
NC_007982.1 RefSeq overlapping_feature_set 182895 182977 + . name_left=gene-ZeamM089;source_left=NC_007982.1|gene|182895|182977|+|gene-ZeamM089;name_right=gene-ZeamM089;source_right=NC_007982.1|gene|182895|182977|+|gene-ZeamM089;name_right=gene-ZeamM089|
NC_007982.1 RefSeq overlapping_feature_set 111207 111278 + . name_left=gene-ZeamM018;source_left=NC_007982.1|gene|111207|111278|+|gene-ZeamM018;name_right=gene-ZeamM018;source_right=NC_007982.1|gene|111207|111278|+|gene-ZeamM018;name_right=gene-ZeamM018|
NC_007982.1 RefSeq overlapping_feature_set 116164 116206 + . name_left=gene-ZeamP035;source_left=NC_007982.1|gene|116164|116206|+|gene-ZeamP035;name_right=gene-ZeamP035;source_right=NC_007982.1|gene|116164|116206|+|gene-ZeamP035;name_right=gene-ZeamP035|
NC_007982.1 RefSeq gene 127146 122538 + . name=gene-ZeamP017;source=NC_007982.1|gene|127146|122538|+|gene-ZeamP017|
NC_007982.1 RefSeq overlapping_feature_set 123748 123822 + . name_left=gene-ZeamM011;source_left=NC_007982.1|gene|123748|123822|+|gene-ZeamM011;name_right=gene-ZeamM011;source_right=NC_007982.1|gene|123748|123822|+|gene-ZeamM011;name_right=gene-ZeamM011|
NC_007982.1 RefSeq overlapping_feature_set 125072 125943 + . name_left=gene-ZeamM012;source_left=NC_007982.1|gene|125072|125943|+|gene-ZeamM012;name_right=gene-ZeamM012;source_right=NC_007982.1|gene|125072|125943|+|gene-ZeamM012;name_right=gene-ZeamM012|
NC_007982.1 RefSeq overlapping_feature_set 128788 128897 + . name_left=gene-ZeamM013;source_left=NC_007982.1|gene|128788|128897|+|gene-ZeamM013;name_right=gene-ZeamM013;source_right=NC_007982.1|gene|128788|128897|+|gene-ZeamM013;name_right=gene-ZeamM013|
NC_007982.1 RefSeq overlapping_feature_set 129585 129658 + . name_left=gene-ZeamM014;source_left=NC_007982.1|gene|129585|129658|+|gene-ZeamM014;name_right=gene-ZeamM014;source_right=NC_007982.1|gene|129585|129658|+|gene-ZeamM014;name_right=gene-ZeamM014|
NC_007982.1 RefSeq overlapping_feature_set 138286 138357 + . name_left=gene-ZeamM015;source_left=NC_007982.1|gene|138286|138357|+|gene-ZeamM015;name_right=gene-ZeamM015;source_right=NC_007982.1|gene|138286|138357|+|gene-ZeamM015;name_right=gene-ZeamM015|
NC_007982.1 RefSeq overlapping_feature_set 132697 132679 + . name_left=gene-ZeamM016;source_left=NC_007982.1|gene|132697|132679|+|gene-ZeamM016;name_right=gene-ZeamM016;source_right=NC_007982.1|gene|132697|132679|+|gene-ZeamM016;name_right=gene-ZeamM016|
NC_007982.1 RefSeq overlapping_feature_set 138489 139212 + . name_left=rna-ZeamM017;source_left=NC_007982.1|rRNA|138489|139212|+|rna-ZeamM017;name_right=gene-ZeamP051;source_right=NC_007982.1|rRNA|138489|139212|+|rna-ZeamM017;name_right=gene-ZeamP051|
NC_007982.1 RefSeq overlapping_feature_set 168734 164661 + . name_left=gene-ZeamP057;source_left=NC_007982.1|gene|168734|164661|+|gene-ZeamP057;name_right=gene-ZeamP024;source_right=NC_007982.1|gene|168734|164661|+|gene-ZeamP057;name_right=gene-ZeamP024|
NC_007982.1 RefSeq rRNA 179618 179691 + . name=rna-ZeamM018;source=NC_007982.1|rRNA|179618|179691|+|rna-ZeamM018|

```

**Figure 3. TIGRE's <>\_clean.gff3 output file. TIGRE's <>\_clean.gff3 file is a conventional gff3 file, but with a tailored attributes column**

The information in the attributes column (i.e., name/name\_left/name\_right and source/source\_left/source\_right) is derived from the elements delimiting each intergenic region. This contextual information is carried over to the subsequent steps, so the user can pinpoint where each intergenic comes from (based not only on genomic coordinates, but also on flanking elements). Detailed explanation of the values from the attributes column can be found in the <https://github.com/brenodupin/tigre/blob/master/PROCESSING.md>.

**Note:** tigre extract was designed to work with the clean GFF3 files generated by tigre clean. It is not certain whether tigre extract will produce proper results with GFF3 files that have not been formatted by tigre clean.

### Running tigre getfasta

⌚ **Timing:** Seconds to minutes (depending on the size of the dataset)

This section fetches the nucleotide sequences of the annotated intergenic regions. This step requires the FASTA files of corresponding genomes.

- Run tigre getfasta to generate FASTA files containing the intergenic regions from <AN>\_intergenic.gff3:

```

>tigre getfasta multiple -v --log tigre_getfasta.log --tsv plants_mit.tsv --bedtools-compatible --overwrite

```

**Note:** tigre getfasta is optional and requires as input the respective FASTA files containing the genomes being analyzed. Although this step is optional, most users might implement it to obtain the intergenic nucleotide sequences for downstream analyses. tigre getfasta relies on the properly annotated <AN>\_intergenic.gff3 files (generated by tigre extract) to handle circular genomes that possess intergenic regions spanning the end/start of the genome.

**CRITICAL:** Even though biopython is implemented to extract the fasta sequences, the flag --bedtools-compatible generates fasta files that are compatible with BEDtools (i.e., the nucleotide sequences will be 0-based). This functionality will not alter the nucleotide sequence itself but only the header information presented in the header (after the ">" sign).

### EXPECTED OUTCOMES

The expected outcomes of this protocol are the annotation and extraction of intergenic regions from any annotated genome. tigre clean will generate, with the optional help of GDT, "clean" GFF3 files

```

igff-version 3
#igff-spec-version 1.26
#processor TIGRE igff.py
#sequence-region NC_007982.1 1 569630
#species https://www.ncbi.nlm.nih.gov/Taxonomy/Browser/wwwtax.cgi?d=381124
NC_007982.1 RefSeq intergenic_region 26890 26899 + ID=NC_007982.1_intergenic_region_2;name_up=gene-ZeamM088;source_up=NC_007982.1|gene|19813|26889|+|gene-ZeamM088;name_dw=gene-ZeamM088;source_dw=NC_007982.1|gene|19813|26889|+|gene-ZeamM088;
NC_007982.1 RefSeq intergenic_region 27213 42662 + ID=NC_007982.1_intergenic_region_3;name_up=gene-ZeamM012;source_up=NC_007982.1|gene|26918|27212|+|gene-ZeamM012;name_dw=gene-ZeamM012;source_dw=NC_007982.1|gene|26918|27212|+|gene-ZeamM012;
NC_007982.1 RefSeq intergenic_region 42725 43261 + ID=NC_007982.1_intergenic_region_4;name_up=gene-ZeamM081;source_up=NC_007982.1|gene|42662|42724|+|gene-ZeamM081;name_dw=gene-ZeamM081;source_dw=NC_007982.1|gene|42662|42724|+|gene-ZeamM081;
NC_007982.1 RefSeq intergenic_region 43436 44162 + ID=NC_007982.1_intergenic_region_5;name_up=gene-ZeamM082;source_up=NC_007982.1|gene|43362|43435|+|gene-ZeamM082;name_dw=gene-ZeamM082;source_dw=NC_007982.1|gene|43362|43435|+|gene-ZeamM082;
NC_007982.1 RefSeq intergenic_region 44237 47076 + ID=NC_007982.1_intergenic_region_6;name_up=gene-ZeamM083;source_up=NC_007982.1|gene|44163|44236|+|gene-ZeamM083;name_dw=gene-ZeamM083;source_dw=NC_007982.1|gene|44163|44236|+|gene-ZeamM083;
NC_007982.1 RefSeq intergenic_region 47189 49197 + ID=NC_007982.1_intergenic_region_7;name_up=gene-ZeamM084;source_up=NC_007982.1|gene|47077|47148|+|gene-ZeamM084;name_dw=gene-ZeamM084;source_dw=NC_007982.1|gene|47077|47148|+|gene-ZeamM084;
NC_007982.1 RefSeq intergenic_region 49273 54809 + ID=NC_007982.1_intergenic_region_8;name_up=gene-ZeamM085;source_up=NC_007982.1|gene|49198|49272|+|gene-ZeamM085;name_dw=gene-ZeamM085;source_dw=NC_007982.1|gene|49198|49272|+|gene-ZeamM085;
NC_007982.1 RefSeq intergenic_region 50875 57193 + ID=NC_007982.1_intergenic_region_9;name_up=gene-ZeamM186;source_up=NC_007982.1|gene|50499|50874|+|gene-ZeamM186;name_dw=gene-ZeamM186;source_dw=NC_007982.1|gene|50499|50874|+|gene-ZeamM186;
NC_007982.1 RefSeq intergenic_region 59341 63967 + ID=NC_007982.1_intergenic_region_10;name_up=gene-ZeamM019;source_up=NC_007982.1|gene|57184|59348|+|gene-ZeamM019;name_dw=gene-ZeamM019;source_dw=NC_007982.1|gene|57184|59348|+|gene-ZeamM019;
NC_007982.1 RefSeq intergenic_region 64862 76131 + ID=NC_007982.1_intergenic_region_11;name_up=gene-ZeamM091;source_up=NC_007982.1|gene|63968|64861|+|gene-ZeamM091;name_dw=gene-ZeamM091;source_dw=NC_007982.1|gene|63968|64861|+|gene-ZeamM091;
NC_007982.1 RefSeq intergenic_region 84977 94266 + ID=NC_007982.1_intergenic_region_12;name_up=gene-ZeamM032;source_up=NC_007982.1|gene|74332|84976|+|gene-ZeamM032;name_dw=gene-ZeamM032;source_dw=NC_007982.1|gene|74332|84976|+|gene-ZeamM032;
NC_007982.1 RefSeq intergenic_region 94338 95773 + ID=NC_007982.1_intergenic_region_13;name_up=gene-ZeamM087;source_up=NC_007982.1|gene|94257|94329|+|gene-ZeamM087;name_dw=gene-ZeamM087;source_dw=NC_007982.1|gene|94257|94329|+|gene-ZeamM087;
NC_007982.1 RefSeq intergenic_region 98074 108984 + ID=NC_007982.1_intergenic_region_14;name_up=gene-ZeamM033;source_up=NC_007982.1|gene|95774|98073|+|gene-ZeamM033;name_dw=gene-ZeamM033;source_dw=NC_007982.1|gene|95774|98073|+|gene-ZeamM033;
NC_007982.1 RefSeq intergenic_region 108979 101543 + ID=NC_007982.1_intergenic_region_15;name_up=rna-ZeamM088;source_up=NC_007982.1|tRNA|108985|108978|+|rna-ZeamM088;name_dw=rna-ZeamM088;source_dw=NC_007982.1|tRNA|108985|108978|+|rna-ZeamM088;
NC_007982.1 RefSeq intergenic_region 102117 102094 + ID=NC_007982.1_intergenic_region_16;name_up=gene-ZeamM036;source_up=NC_007982.1|gene|102544|102116|+|gene-ZeamM036;name_dw=gene-ZeamM036;source_dw=NC_007982.1|gene|102544|102116|+|gene-ZeamM036;
NC_007982.1 RefSeq intergenic_region 102978 111286 + ID=NC_007982.1_intergenic_region_17;name_up=gene-ZeamM089;source_up=NC_007982.1|gene|102895|102977|+|gene-ZeamM089;name_dw=gene-ZeamM089;source_dw=NC_007982.1|gene|102895|102977|+|gene-ZeamM089;
NC_007982.1 RefSeq intergenic_region 111279 111613 + ID=NC_007982.1_intergenic_region_18;name_up=gene-ZeamM018;source_up=NC_007982.1|gene|111287|111278|+|gene-ZeamM018;name_dw=gene-ZeamM018;source_dw=NC_007982.1|gene|111287|111278|+|gene-ZeamM018;
NC_007982.1 RefSeq intergenic_region 115827 122145 + ID=NC_007982.1_intergenic_region_19;name_up=gene-ZeamM015;source_up=NC_007982.1|gene|111614|115826|+|gene-ZeamM015;name_dw=gene-ZeamM015;source_dw=NC_007982.1|gene|111614|115826|+|gene-ZeamM015;
NC_007982.1 RefSeq intergenic_region 122531 122747 + ID=NC_007982.1_intergenic_region_20;name_up=gene-ZeamM012;source_up=NC_007982.1|gene|122146|122530|+|gene-ZeamM012;name_dw=gene-ZeamM012;source_dw=NC_007982.1|gene|122146|122530|+|gene-ZeamM012;
NC_007982.1 RefSeq intergenic_region 123823 125871 + ID=NC_007982.1_intergenic_region_21;name_up=gene-ZeamM011;source_up=NC_007982.1|gene|123748|123822|+|gene-ZeamM011;name_dw=gene-ZeamM011;source_dw=NC_007982.1|gene|123748|123822|+|gene-ZeamM011;
NC_007982.1 RefSeq intergenic_region 125844 128783 + ID=NC_007982.1_intergenic_region_22;name_up=gene-ZeamM012;source_up=NC_007982.1|gene|125872|125843|+|gene-ZeamM012;name_dw=gene-ZeamM012;source_dw=NC_007982.1|gene|125872|125843|+|gene-ZeamM012;

```

**Figure 4. TIGRE's <>\_intergenic.gff3 output file. TIGRE's <>\_intergenic.gff3 file is a conventional GFF3 file containing only the intergenic regions of the inputted genome**

Note that, just like in. <>\_clean.gff3, the attributes column of <>\_intergenic.gff3 is customized and carries the contextual information necessary for localizing each intergenic region. <https://github.com/brenodupin/tigre/blob/master/PROCESSING.md> contains further details on how TIGRE deals with elements (an intergenic region or otherwise) spanning the boundary of circular genomes.

formatted for the subsequent commands (tigre extract and tigre getfasta). The clean GFF3 files serve as templates to reveal where the intergenic regions are situated, including the complicated case of intergenic regions wrapping around the end/start boundaries of circular genomes. tigre extract generates “<>\_intergenic.gff3” files that have the intergenic regions annotated and ready to be extracted by tigre getfasta. Via biopython, tigre getfasta can extract hundreds of thousands of intergenic sequences in a matter of seconds. With its user-centered design along with its multiprocessing capabilities, TIGRE is meant to facilitate large-scale comparative analyses of intergenic regions across disparate species.

One extra outcome of this protocol is the possible advanced usage of TIGRE to extract intronic regions. As introns represent another major component in the noncoding portions of genomes, the command tigre combine has been developed to be used as a workaround to extract introns. In a staggered approach via tigre clean and tigre combine, the user can create a “clean\_to\_introns” GFF3 file that will serve as input for the extraction of introns via tigre extract. See [troubleshooting](https://github.com/brenodupin/tigre) (and Advanced Usage in <https://github.com/brenodupin/tigre>) for more details. Note that this usage is experimental, as it is a process for the extraction (and not the identification) of introns from exonic annotations.

## LIMITATIONS

TIGRE has two main limitations. The first is that this protocol was designed with GFF3 files that closely resemble the GFF3's from NCBI RefSeq. Essential parts of TIGRE rely on the existence (or absence) of elements common to these files, such as “region” being the feature type of the first line after the header, the presence of “is\_circular=true” to confirm the circularity of a genome, and the field “ID=” in the column attributes to provide the names of the genes that delimit intergenic regions (regardless whether GDT is implemented or not). One could argue that TIGRE should have been designed with other annotation files in mind, such as BED and GTF, but these files do not have all the necessary information for the type of downstream analyses that TIGRE permits. GTF's (which are identical to GFF2 - <https://gmod.org/wiki/GFF2>) contain hierarchical information up to only two levels (exons to transcripts and transcripts to genes). Meanwhile, intervening sequences can be easily extract from BED files with BEDtools via bedtools complement (<https://bedtools.readthedocs.io/en/latest/content/tools/complement.html>). But again, no identifying information of the delimiting features can be extracted from a BED file.

Regarding file structure, \_intergenic.gff3 files are not readily “mergeable” to a conventional GFF3 file. This is because tigre extract outputs GFF3 files that follow the format created by tigre clean. Particularly, the clean GFF3 attributes column is very different than the attributes column found in a conventional GFF3 (from having information such as name\_up/source\_up and name\_dw/source\_dw to having every line ending in “;”, which is not always the case for regular GFF3 files).

## TROUBLESHOOTING

### Problem 1

What if I want to extract introns as well?

### Potential solution

With the `_intergenic.gff3` in hand, merge the `_intergenic.gff3` with the original GFF using the `tigre combine` command. In the combined file(s), rerun `tigre clean` with the following modified flags.

```
--query-string "type in ('exon', 'intergenic region', 'tRNA', 'rRNA')"
```

```
--gff-in-suffix "_combined"
```

```
--gff-out-suffix "_clean_to_introns"
```

```
--extended-filtering
```

This will create a clean GFF3 suitable for the extraction of introns only (via `tigre extract` and, possibly, `tigre getfasta`).

**Note:** `--gff-in-suffix` as `"_combined"` refers to the default output of `tigre combine`. `--gff-out-suffix` as `"_clean_to_introns"` is a recommendation for the output suffix, just so the output of this `tigre clean` iteration does not get mixed up with previous outputs of `tigre clean`. Note that `"_clean_to_introns"` must be passed as `--gff-in-suffix` for `tigre extract` (to extract the intronic regions). More details of this implementation can be found at <https://github.com/brenodupin/tigre> under Advanced Usage.

### Problem 2

I followed this protocol using the example dataset (`plants_mit`). How can I compare my results with the results stored in the expected outcome folder (<https://github.com/brenodupin/tigre/tree/master/examples>)?

### Potential solution

You should download `plants_mit_expected_outcome.tar.zst` (preferably in the same directory where the `plants_mit` folder is stored). Then, you should decompress the expected outcome folder with `>tar -xf plants_mit_expected_outcome.tar.zst`. After that, you can compare your results with the expected outcomes by running `>diff -rq plants_mit/plants_mit_expected_outcome/`. If your results are identical to the expected outcome, the only differences should appear in the log files (as they are time stamped).

## RESOURCE AVAILABILITY

### Lead contact

Further information and requests for resources should be directed to and will be fulfilled by the lead contact, Matheus Sanita Lima ([msanital@uwo.ca](mailto:msanital@uwo.ca)).

### Technical contact

Technical questions on executing this protocol should be directed to and will be answered by the technical contact, Breno Dupin ([breno.dupin@gmail.com](mailto:breno.dupin@gmail.com)).

### Materials availability

Not applicable.

### Data and code availability

The datasets and code generated during this study are available at GitHub: <https://github.com/brenodupin/tigre> and Zenodo: <https://zenodo.org/records/17475986>.

## ACKNOWLEDGMENTS

We would like to kindly thank Dr. Igor Fesenko for swiftly responding to our inquiries when we started investigating organelle intergenic sequences. TIGRE started off as an adaptation of Dr. Fesenko's protocol for the extraction of bacterial intergenic sequences. We would also like to thank the Natural Sciences and Engineering Research Council of Canada (NSERC) and the Fundação Araucária and CNPq for financial support granted to D.R.S. (via a NSERC Discovery Grant) and to A.R.P. (via NAPI Bioinformática 66.2021 and 440412/2022-6; Fundação Araucária and CNPq, respectively).

## AUTHOR CONTRIBUTIONS

Conceptualization, M.S.L., B.D., D.R.S., and A.R.P.; software and methodology, B.D.; data curation, M.S.L. and B.D.; writing – original draft, M.S.L. and B.D.; writing – review and editing, M.S.L., B.D., D.R.S., and A.R.P.; funding acquisition and supervision, D.R.S. and A.R.P.

## DECLARATION OF INTERESTS

The authors declare no competing interests.

## DECLARATION OF GENERATIVE AI AND AI-ASSISTED TECHNOLOGIES IN THE WRITING PROCESS

During the preparation of this work, the authors used Claude AI and ChatGPT for script optimization and documentation, code formatting, GitHub backend configuration, tool logo generation, and alike. After using this tool/service, the authors reviewed and edited the content as needed and take full responsibility for the content of the published article.

## REFERENCES

- Eddy, S.R. (2013). The ENCODE project: Missteps overshadowing a success. *Curr. Biol.* 23, R259–R261.
- Graur, D., Zheng, Y., Price, N., Azevedo, R.B.R., Zufall, R.A., and Elhaik, E. (2013). On the immortality of television sets: “function” in the human genome according to the evolution-free gospel of ENCODE. *Genome Biol. Evol.* 5, 578–590.
- Doolittle, W.F. (2013). Is junk DNA bunk? A critique of ENCODE. *Proc. Natl. Acad. Sci. USA* 110, 5294–5300.
- Lynch, M., and Conery, J.S. (2003). The origins of genome complexity. *Science* 302, 1401–1404.
- Xu, H., Li, C., Xu, C., and Zhang, J. (2023). Chance promoter activities illuminate the origins of eukaryotic intergenic transcriptions. *Nat. Commun.* 14, 1826.
- Sanita Lima, M., and Smith, D.R. (2017). Pervasive transcription of mitochondrial, plastid, and nucleomorph genomes across diverse plastid-bearing species. *Genome Biol. Evol.* 9, 2650–2657.
- Werner, A., Kanhere, A., Wahlestedt, C., and Mattick, J.S. (2024). Natural antisense transcripts as versatile regulators of gene expression. *Nat. Rev. Genet.* 25, 730–744.
- Attallah, C., Conti, G., Zuljan, F., Zavallo, D., and Ariel, F. (2025). Noncoding RNAs as tools for advancing translational biology in plants. *Plant Cell* 37, koaf054.
- Narunsky, A., Higgs, G.A., Torres, B.M., Yu, D., de Andrade, G.B., Kavita, K., and Breaker, R.R. (2024). The discovery of novel noncoding RNAs in 50 bacterial genomes. *Nucleic Acids Res.* 52, 5152–5165.
- Lynch, M. (2006). Streamlining and simplification of microbial genome architecture. *Annu. Rev. Microbiol.* 60, 327–349.
- Wu, P., Xue, N., Yang, J., Zhang, Q., Sun, Y., and Zhang, W. (2025). OGU: a toolbox for better utilising organelle genomic data. *Mol. Ecol. Resour.* 25, e14044.
- Fesenko, I., Sahakyan, H., Dhyani, R., Shabalina, S.A., Storz, G., and Koonin, E.V. (2025). The hidden bacterial microproteome. *Mol. Cell* 85, 1024–1041.e6.
- Thorpe, H.A., Bayliss, S.C., Hurst, L.D., and Feil, E.J. (2017). Comparative analyses of selection operating on nontranslated intergenic regions of diverse bacterial species. *Genetics* 206, 363–376.
- Tsai, C.-H., Liao, R., Chou, B., Palumbo, M., and Contreras, L.M. (2015). Genome-wide analyses in bacteria show small-RNA enrichment for long and conserved intergenic regions. *J. Bacteriol.* 197, 40–50.
- Dupin, B., Sanita Lima, M., Rossi Paschoal, A., and Smith, D.R. (2025). Protocol for GDT, Gene Dictionary Tool, to create and implement a gene dictionary across annotated genomes. *STAR Protocols* 6, 104187. <https://doi.org/10.1016/j.xpro.2025.104187>.
- Cock, P.J.A., Antao, T., Chang, J.T., Chapman, B.A., Cox, C.J., Dalke, A., Friedberg, I., Hamelryck, T., Kauff, F., Wilczynski, B., and de Hoon, M.J.L. (2009). Biopython: freely available Python tools for computational molecular biology and bioinformatics. *Bioinformatics* 25, 1422–1423.