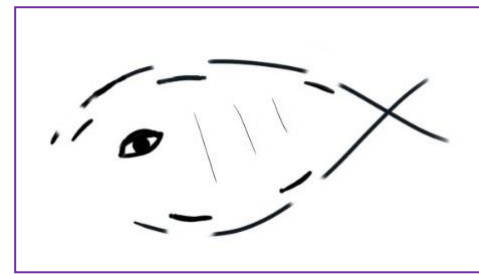


NPGREAT

NanoPore Guided REgional Assembly Tool



1. Method

NPGREAT is a hybrid assembly method that utilizes two complementary types of data for the assembly of the human subtelomere regions: The Linked-Reads & the Nanopore reads.

2. Requirements

- Programs: Blastn, Repeat Masker, Tandem Repeat Finder, Seqtk
- Python libraries: Biopython, Pandas

3. Data

All data correspond to the NA12878 human cell line.

- REXTAL assembly contigs obtained by assembling 10x Linked-Reads
- REXTAL assembly contigs obtained by assembling Tell-Seq Linked-Reads
- Telomeric and Subtelomeric ONT reads
- NPGREAT output assembly using as input the 10x and ONT data
- NPGREAT output assembly using as input the Tell-Seq and ONT data

4. Code

The code is in python scripts and a shell script.

The code is divided in the NPGREAT steps: npgreat_1orientation.sh, npgreat_2position.py, npgreat_3correction.py, npgreat_4connectors_gapfilling_combination.py

5. Usage

- **nano_telom_file:** The FASTA file containing the telomeric Nanopore reads^{1,2}.
- **nano_subtelom_file:** The FASTA file containing the subtelomeric Nanopore reads¹.
- **rextal_contigs_file:** The FASTA file containing the REXTAL assembly contigs³.
- **subtel_region_name:** The name of the subtelomeric region to assemble in the format <chromosome><arm>.
- **repeat_masker_exec:** The path to the Repeat Masker executable.
- **tandem_repeat_finder_exec:** The path to the Tandem Repeat Finder executable.
- **eval_option:** The evalute to be used in the blastn alignments during the steps. Select 0 for the zero blastn evalute (recommended for most subtelomeric regions) or set to 1 for the default blastn evalute (recommended for subtelomeric regions 16p, 19q, 20p & 22q).

```
./npgreat_1orientation.sh [nano_telom_file] [nano_subtelom_file] [rextal_contigs_file]  
[subtel_region_name] [repeat_masker_exec] [tandem_repeat_finder_exec]
```

```
python npgreat_2position.py [eval_option]
```

```
python npgreat_3correction.py [eval_option]
```

```
python npgreat_4connectors_gapfilling_combination.py
```

¹. The ultralong Nanopore reads that are above 40Kbp. ². The Nanopore reads that contain the telomere repeat tract TTAGGGx4. ³. The output FASTA file of the REXTAL method after: (i) the 10Ns & 100Ns have been removed from the scaffolds, (ii) the scaffolds have been split at the location where 100Ns existed, (iii) all the newly created contigs of the file have unique IDs.

6. Example of running NPGREAT

An example of using NPGREAT to compute the assembly of the 10p subtelomeric region. The input data are the (sub)telomeric Nanopore reads and the Tell-Seq Linked-Read REXTAL assembly contigs.

➤ We run the shell script:

```
[npgreat_folder]$ ./npgreat_1orientation.sh telom.fasta subtelom.fasta rextal_contigs.fasta 10p /home/programs/RepeatMasker /home/programs/trf/trf409.legacylinux64

*** The orientation step begins... ***

RepeatMasker output...
analyzing file others.fa
...

RepeatMasker output...
analyzing file oriented_nanos.fa
...

*** The nanopore reads have been oriented. ***
*** The REXTAL contigs have been oriented. ***
*** The orientation step finished. ***
*****
```

➤ We run the python files⁴:

```
[npgreat_folder]$ python npgreat_2position.py 0

- The NPGREAT Position step begins...
Calculating the alignments...
Extracting position information...
The NPGREAT Position step finished.

[npgreat_folder]$ python npgreat_3correction.py 0

- The NPGREAT Correction step begins...
Detecting possible splits by checking the internal alignments...
Investigating detected possible splits...
Identifying splits...
Updating position information of splitted contigs...
The NPGREAT Correction step finished.

[npgreat_folder]$ python npgreat_4connectors_gapfilling_combination.py

- The NPGREAT Region Extraction/Connector Segments step begins...
The NPGREAT Region Extraction/Connector Segments step finished.
- The NPGREAT Gap Filling step begins...
The NPGREAT Gap Filling step finished.
- The NPGREAT Combination step begins...
The NPGREAT Combination step finished.
*****
*** The NPGREAT assembly has been computed. Total assembly length: 410997bp ***
Please see output file: assembly_npgreat.fasta
*****
```

⁴. We used a terminal to run the python files, but the code can also run from a Python program or IDE, like for example Spider or Jupyter Notebook.

