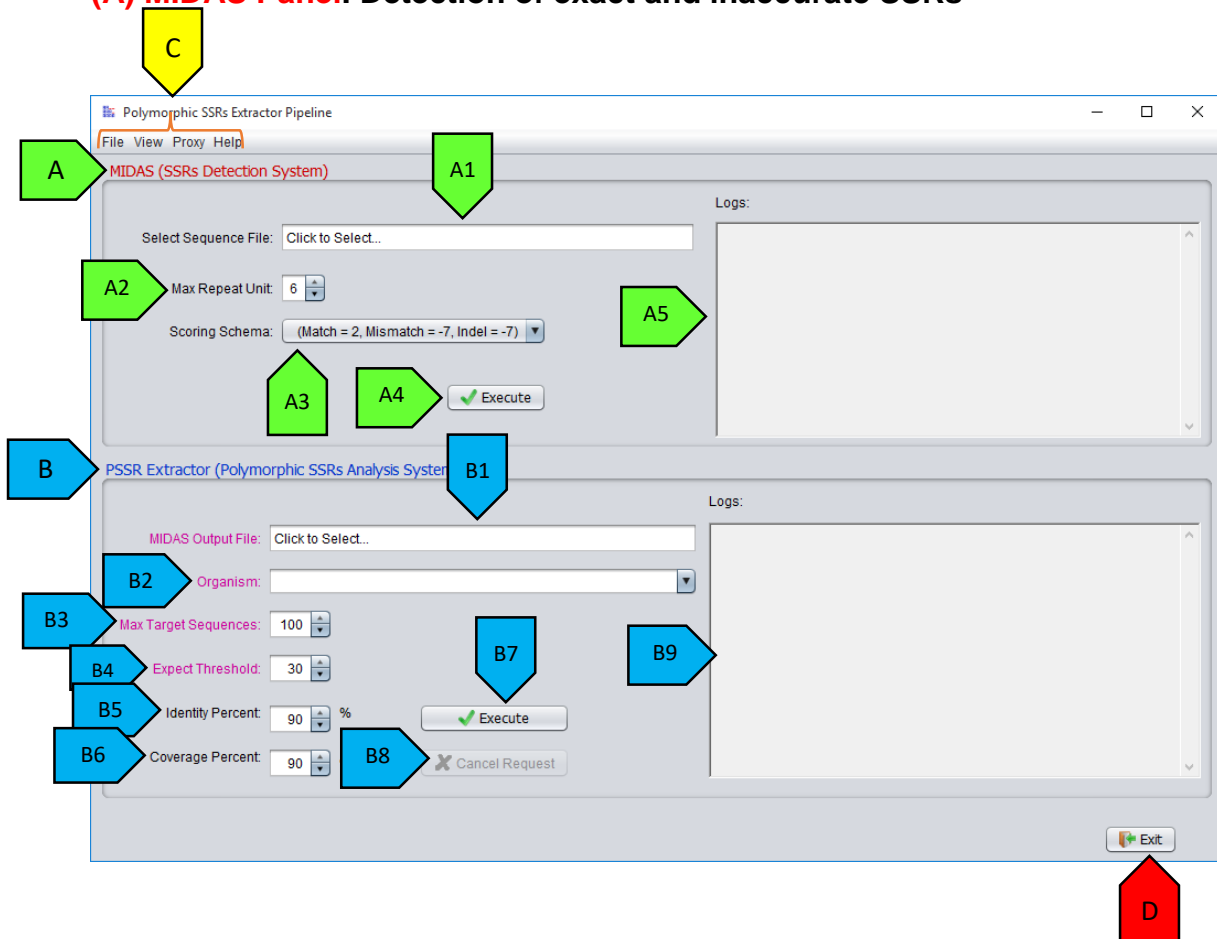


Background

Polymorphic SSRs Extractor Pipeline is a software that allows the fast and reliable identification of polymorphic SSRs loci in genomic sequences. Global processing is done by concatenation of MIDAS [1] and PSSR-Extractor [2] programs, the latter making a remote invocation to NCBI BLAST [3].

(A) MIDAS Panel. Detection of exact and inaccurate SSRs



(A1) Select a file (genome in FASTA or GBFF format, both single or multi-locus).

(A2) Sets the maximum size of the repeat unit to be scanned in the genome (1-6).

(A3) Alignment parameters scheme for match, mismatch and indel.

(A4) Execute the program according to the previously defined parameters.

(A5) Reports the current status of the application execution.

.xls file:

Accession: GenBank access code.

Length: Size of the repeating unit.

Start: Start position in the genome.

End: End position in the genome.

Score: Alignment score.

Matches: Matching bases.

Mismatches: Mismatched bases.

Indel: Base insertions and deletions.

Inaccuracy: SSR inaccuracy (%).

L_Flank: Left flank sequence.

Entropy_LF: Compositional entropy in left flank.

R_Flank: Right flank sequence.

Entropy_RF: Compositional entropy in right flank.

.dat file. It could be visualized in any text editor, presents the data in non-tabular form and shows sequence alignments:

```

Assembly: G:\generos\Bacterias_2019\GCF_000008585.1_ASM858v1_genomic.gbff
Max Pattern Length= 6
Alignment's Parameters: Match=2 Mismatch=-7 Indel=-7
-----
Accession: NC_002755.2
Found Repeats= 198
-----
Pattern:                aatacgc
Exact SSR Pos.:         1112963
Exact SSR Copy Number:  4
Matches:                26
Mismatches:             0
Insertions/Deletions:   0
Inexact SSR Start Pos.: 1112962
Inexact SSR End Pos.:   1112987
Inexact SSR Copy Number: 4
Score:                  52
L_Flank:  gacccggaggccgacccggt  Entropy(0-2): 1.68614
R_Flank:  tcgaggacacctgcggtttg  Entropy(0-2): 1.94188
Aligned SSR:
cgaatacgaatacgaatacgaatacgc
|||||
cgaatacgaatacgaatacgaatacgc
-----

```

.mfaa file. This file is the input for PSSR-Extractor program (B). It could be visualized in any text editor and presents the SSRs in multi-fasta format. The repeat region is marked in lowercase and the flanks in uppercase. The header presents information such as the GenBank accession number, the motif, and the positions in the genome:

```

>NC_002755.2|aatacgc[1112962-1112987]|c:4|s:52|m:26|mm:0|i:0|ina:0
|5f:gacccggaggccgacccggt|5e:1.69|3f:tcgaggacacctgcggtttg|3e:1.94
GACCCGGAGGCCGACCCGGTcgaatacgaatacgaatacgaatacgcTCGAGGACACCTGCGGTTT
G
>NC_002755.2|accagc[2413197-2413217]|c:3|s:42|m:21|mm:0|i:0|ina:0
|5f:ggcctccttgccgatccccg|5e:1.68|3f:tcacgatggtgatcgcgaaa|3e:1.97
GGCCTCCTTGCCGATCCCCGgcaccagcaccagcaccagcaTCACGATGGTGATCGCGAAA
>NC_002755.2|accgcc[3869815-3869832]|c:3|s:36|m:18|mm:0|i:0|ina:0
|5f:aaccgctagccccacagttg|5e:1.91|3f:caacgccaggccctgatcg|3e:1.86
AACCGCTAGCCCCACAGTTGacggccaccgccaccgccCAACGCCAGGGCCTGATCGG
>NC_002755.2|acggcg[1339305-1339330]|c:4|s:52|m:26|mm:0|i:0|ina:0
|5f:cgcgatgtttggctacgccc|5e:1.88|3f:ttgctgcggttcgaggaggc|3e:1.86
CGCGATGTTTGGCTACGCCGcgggcgacggcgacggcgacggcgacgTTGCTGCCGTTCGAGGAGG
C
>NC_002755.2|atgtcg[3775640-3775660]|c:3|s:42|m:21|mm:0|i:0|ina:0
|5f:ccttttcgcgctgatccgac|5e:1.85|3f:ggtggcccgctcgcgggcg|3e:1.3
CCTTTTCGCGCTGATCCGACTcgatgtcgatgtcgatgtcgGGTGGCCCGTCGCGGGCGGG
>NC_002755.2|cccgcg[3802087-3802105]|c:3|s:38|m:19|mm:0|i:0|ina:0
|5f:gaaggccacgggccaccact|5e:1.74|3f:tcaccggcgccctccaga|3e:1.68
GACGGGCACGGGCCACCACTgcccgcgcccgcgcccgcgTCACCGGCGCCCTCCAGAA
>NC_002755.2|accgg[3107247-3107262]|c:3|s:32|m:16|mm:0|i:0|ina:0|
5f:tcggcgccgggcccggggcca|5e:1.46|3f:tcacagtgcctcctcgct|3e:1.72
TCGCGGCCGGGCCGGGGCCAaccggacccgacccgaTCCAGTGCTCGTCCCTCGCT

```

(B) PSSR Extractor Panel. Detects polymorphism in microsatellites based on MIDAS results.



(B1) Selects the .mfaa file from MIDAS, which will be the queries for BLAST.

(B2) Selects the organism to which the sequences belong. BLAST will search sequences from the same organism.

(B3) Sets the maximum of similar sequences that BLAST will return.

(B4) Filters sequences with an expected value higher than the selected one.

(B5) Minimum percent identity that the two sequences must have on their flanks.

(B6) Percentage of the number of bases that the two sequences must share on their flanks.

(B7) Execute the program according to the previously defined parameters.

(B8) Cancel the execution of the program. This button is activated after pressing B7.

(B9) Reports the current state of application execution (B).

(D) Exits the program.

BLAST has more parameters than those exposed, but it is necessary for the analysis kept them constant. These parameters and their default value are shown below:

QUERY_BELIEVE_DEFLINE: false.

DATABASE: nr.

LCASE_MASK: true.

FILTER: F.

FORMAT_TYPE: Tabular.

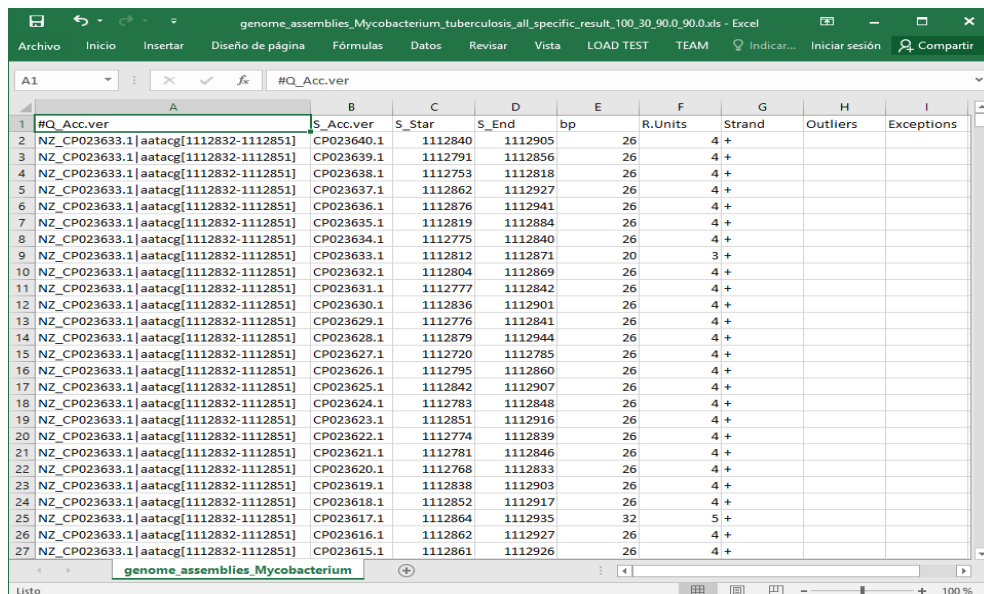
PROGRAM: blastn.

CLIENT: web.

BLAST_PROGRAM: blastn.

PSSR Extractor outputs two .xls files, one detailed and the other generic, whose names have the suffixes **_specific_result** and **_generic_result** respectively.

_specific_result.xls. Provides detailed information on each processed subject.



#Q_Acc.ver	S_Acc.ver	S_Start	S_End	bp	R.Units	Strand	Outliers	Exceptions
NZ_CP023633.1 aatagc[1112832-1112851]	CP023640.1	1112840	1112905	26	4 +			
NZ_CP023633.1 aatagc[1112832-1112851]	CP023639.1	1112791	1112856	26	4 +			
NZ_CP023633.1 aatagc[1112832-1112851]	CP023638.1	1112753	1112818	26	4 +			
NZ_CP023633.1 aatagc[1112832-1112851]	CP023637.1	1112862	1112927	26	4 +			
NZ_CP023633.1 aatagc[1112832-1112851]	CP023636.1	1112876	1112941	26	4 +			
NZ_CP023633.1 aatagc[1112832-1112851]	CP023635.1	1112819	1112884	26	4 +			
NZ_CP023633.1 aatagc[1112832-1112851]	CP023634.1	1112775	1112840	26	4 +			
NZ_CP023633.1 aatagc[1112832-1112851]	CP023633.1	1112812	1112871	20	3 +			
NZ_CP023633.1 aatagc[1112832-1112851]	CP023632.1	1112804	1112869	26	4 +			
NZ_CP023633.1 aatagc[1112832-1112851]	CP023631.1	1112777	1112842	26	4 +			
NZ_CP023633.1 aatagc[1112832-1112851]	CP023630.1	1112836	1112901	26	4 +			
NZ_CP023633.1 aatagc[1112832-1112851]	CP023629.1	1112776	1112841	26	4 +			
NZ_CP023633.1 aatagc[1112832-1112851]	CP023628.1	1112879	1112944	26	4 +			
NZ_CP023633.1 aatagc[1112832-1112851]	CP023627.1	1112720	1112785	26	4 +			
NZ_CP023633.1 aatagc[1112832-1112851]	CP023626.1	1112795	1112860	26	4 +			
NZ_CP023633.1 aatagc[1112832-1112851]	CP023625.1	1112842	1112907	26	4 +			
NZ_CP023633.1 aatagc[1112832-1112851]	CP023624.1	1112783	1112848	26	4 +			
NZ_CP023633.1 aatagc[1112832-1112851]	CP023623.1	1112851	1112916	26	4 +			
NZ_CP023633.1 aatagc[1112832-1112851]	CP023622.1	1112774	1112839	26	4 +			
NZ_CP023633.1 aatagc[1112832-1112851]	CP023621.1	1112781	1112846	26	4 +			
NZ_CP023633.1 aatagc[1112832-1112851]	CP023620.1	1112768	1112833	26	4 +			
NZ_CP023633.1 aatagc[1112832-1112851]	CP023619.1	1112838	1112903	26	4 +			
NZ_CP023633.1 aatagc[1112832-1112851]	CP023618.1	1112852	1112917	26	4 +			
NZ_CP023633.1 aatagc[1112832-1112851]	CP023617.1	1112864	1112935	32	5 +			
NZ_CP023633.1 aatagc[1112832-1112851]	CP023616.1	1112862	1112927	26	4 +			
NZ_CP023633.1 aatagc[1112832-1112851]	CP023615.1	1112861	1112926	26	4 +			

Q_acc_ver: Access version of the query sequence that was compared with the subject sequences in database.

S_acc_ver: Access version of the subject sequences in database that was compared with query sequence.

S_start: Start position of the subject.

S_end: End position of the subject.

Note: The positions provided by S_start and S_end include the flanks surrounding the microsatellite. If S_start > S_end the subject sequence has a negative direction ("-"), that is, a reverse complement sequence.

bp: Number of nucleotides between two flanks.

R.units: Number of units repeated between two flanks.

Outliers: Represented with "***" and can vary from 1 to 5. It means that the number of repeated units between the flanks of the subjects is doubtful because it is very large, being unlikely a microsatellite between them. The cutoff values established for this exception were mononucleotide: 157 bp, dinucleotide: 364 bp, trinucleotide: 109 bp, tetranucleotide: 45 bp, pentanucleotide: 150 bp, and hexanucleotide: 193 bp. These values were defined after processing SSRs in 200 bacterial genomes, recording their sizes, and establishing the cutoff at 3 times the interquartile range:

- * **bp** > cutoff value.
- ** **bp** ≥ twice the cutoff value.
- *** **bp** ≥ three times the cutoff value.
- **** **bp** ≥ four times the cutoff value.
- ***** **bp** ≥ five times the cutoff value.

"D" (degenerated). Subjects have an Identity Percent <90% and / or a Coverage Percent <90%.

"O" (outlier).

"U" (unpair). A single edge appears for the same subject sequence.

Archivo		Inicio	Insertar	Diseño de página	Fórmulas	Datos	Revisar	Vista	LOAD TEST	TEAM	¿Qué desea hacer?	Iniciar sesión	Compartir								
A1																					
A																					
1	MSR_id	Q_Acc.ver	Q_Start	Q_End	Pattern	P_Length	RN	Inaccuracy(N)	L_Flank	Entropy_LF	R_Flank	Entropy_RF	min_RN	max_RN	N	O	Frequency	Q	R	S	Exceptions
2	NZ_CP023633.1 aatacc 1112832-1112851	NZ_CP023633.1	1112832	1112851	aatacc	6	3	0	gaaccgacggccagccggt	1.69	cgatcgacacccctggtttg	1.94	3	5	2		0.02	3	0	0.077	
3	NC_002755.2 accag 2413194-2413220	NC_002755.2	2413194	2413220	accagc	6	4	7.41	gtttccctgcctgcagatc	1.77	cgatcgatcgacggaacac	1.94	4	5	1		0.01	2	0.02		
4	NC_002755.2 accag 3869796-3869840	NC_002755.2	3869796	3869840	accagc	6	7	20	cgctgagacatccagagaa	1.93	ggatcctgacgcgagagcgc	1.69	7	8	1		0.005	2	0.01		
5	NC_002755.2 accag 1339295-1339330	NC_002755.2	1339295	1339330	accagc	6	6	8.11	aaagacgcgcgcgagatttg	1.97	tgctgtgcctgcggagagac	1.86	6	6	0		1	1	0		
6	NC_002755.2 accag 3775640-3775660	NC_002755.2	3775640	3775660	accagc	6	3	0	cccttcctgcctgcctgac	1.85	ggatcctgcctgcctgac	1.53	3	4	1		0.01	2	0.02		
7	NC_002755.2 ccgg 3802087-3802118	NC_002755.2	3802087	3802118	ccggcc	6	5	147	gaacgcgcgcgcgcacacat	1.74	tcctccagcgcgcgcacacac	1.65	5	5	0		1	1	0		
8	NC_002755.2 accag 3107244-3107265	NC_002755.2	3107244	3107265	accagc	6	4	9.09	ccatcgacgcctgcctgcac	1.46	cgatgcctgcctgcctgcac	1.75	4	4	0		1	1	0		
9	NC_002755.2 ccgg 3801343-3801374	NC_002755.2	3801343	3801374	ccggcc	6	5	15.2	ccacacagatgcctgcaggt	1.93	cgatcgatgcctgcctgcaggt	1.92	6	7	1		0.01	2	0.02		
10	NC_002755.2 cgag 1188119-1188140	NC_002755.2	1188119	1188140	cgagcc	5	4	4.55	cgccgcgcgcgcgcgcgag	1.69	tgctgcgcgcgcgcgcgcac	1.86	4	4	0		1	1	0		
11	NC_002755.2 accag 3457937-3457977	NC_002755.2	3457937	3457977	cgagcc	6	8	11.9	gaacgcgcgcgcgcgcgag	1.74	tgctgcgcgcgcgcgcgag	1.91	7	8	1		0.02	2	0.32		
12	NC_002755.2 accag 2560350-2560361	NC_002755.2	2560350	2560361	accagc	4	3	0	gaacgcgcgcgcgcgcctgc	1.93	gattctgcgcgcgcgcgcctgc	1.87	3	3	0		1	1	0		
13	NC_002755.2 accag 1921162-1921173	NC_002755.2	1921162	1921173	accagc	4	3	0	cgatgtgacgcgcgcgcacat	1.94	gcataagacgcgcgcgcgcgc	1.88	3	3	0		1	1	0		
14	NC_002755.2 accag 3198386-3198397	NC_002755.2	3198386	3198397	accagc	4	3	0	ggatggacgcgcgcgcgcgcgc	1.79	cgatcgcgcgcgcgcgcgcgc	1.99	3	3	0		1	1	0		
15	NC_002755.2 atcc 2681722-2681734	NC_002755.2	2681722	2681734	atccgc	4	3	0	gcgcgcgcgcgcgcgcgcgcgc	1.74	gattctgtgcgcgcgcgcgcgc	1.71	3	3	0		1	1	0		
16	NC_002755.2 atgac 3576878-3576899	NC_002755.2	3576878	3576899	atgacg	4	5	9.09	cttatcgacacccgcgcgcgcgc	1.74	ggatgcgcgcgcgcgcgcgcgc	1.41	5	6	1		0.01	2	0.02		
17	NC_002755.2 atgac 1775461-1775473	NC_002755.2	1775461	1775473	atgacg	4	3	0	tcacgcgcgcgcgcgcgcgcgc	1.86	tcctccgcgcgcgcgcgcgcgc	1.95	3	3	0		1	1	0		
18	NC_002755.2 ccgg 2309541-2309552	NC_002755.2	2309541	2309552	ccggcc	4	3	0	ccctgcgcgcgcgcgcgcgcgc	1.93	tcacgcgcgcgcgcgcgcgcgc	1.93	3	3	0		1	1	0		

SSr_id: GenBank query access code, repeating unit and positions in the genome.

Q_acc.ver: GenBank query access code.

Q_Start: Start position of query in genome.

Q_End: End position of query in genome.

Pattern: Pattern sequence.

P_length: Size of the pattern.

RN: Number of pattern sequences.

Inaccuracy: Inaccuracy of the repeat (%).

L_Flank: Flank sequence at left.

Entropy_LF: Compositional entropy in left flank.

R_Flank: Flank sequence at right.

Entropy_RF: Compositional entropy in right flank.

min_RN: Minimum **RN** in all subjects.

max_RN: Maximum **RN** in all subjects.

Range: Difference between **max_RN** and **min_RN**.

Frequency: Allelic frequency of SSR in query.

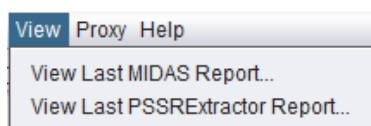
Alleles: Number of alleles.

PIC: Polymorphic Information Content ($1 - \sum_i p_i^2$), also known as expected average of heterozygosity or Nei genetic diversity. It gives a measure of the probability that a pair of randomly chosen alleles in individuals from the same population are different.

Exceptions: Shows all the labels that correspond to exceptions (in **_specific_result.xls**). If all the subject sequences present exceptions then the labels are placed, otherwise the cell appears in blank.

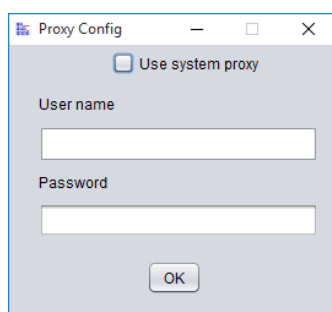
Other options in menu bar.

Report Options in View Menu:



Direct access to the latest MIDAS report.

Direct access to the latest MIDAS report.



It must be checked in case the internet connection is through a proxy. This option uses the proxy settings established by the operating system. If the internet connection through the proxy requires username and password, these fields must be filled in. The "OK" button saves the configuration.

Referencias:

1. Ortíz CMM. MIDAS: Computer application for the identification of exact and inaccurate microsatellites in genomic sequences. Revista Cubana de Informática Médica. 2018;18.
2. Ortíz CMM, Bandinez AR. Methodology for in silico mining of microsatellite polymorphic loci. Revista Cubana de Informática Médica. 2019;19.
3. BLAST Homepage and Selected Search Pages: Introducing the BLAST homepage and form elements/functions of selected search pages2016. Available from: <https://blast.ncbi.nlm.nih.gov>.
4. Fassler J, Cooper P. BLAST Glossary. Available from: <https://blast.ncbi.nlm.nih.gov>.
5. QBLAST's URL API User Guide. Available from: <http://www.ncbi.nlm.nih.gov/blast/Doc/urlapi.html>.