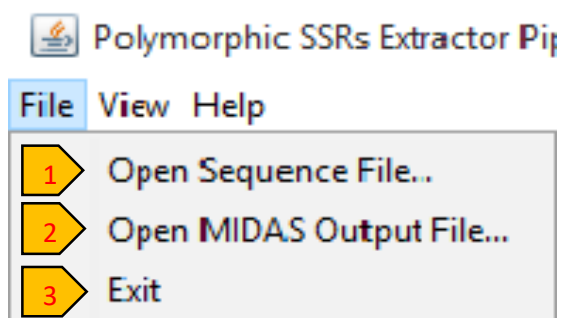


Resumen:

(1) Barra de menú del programa:



(A) Archivo:



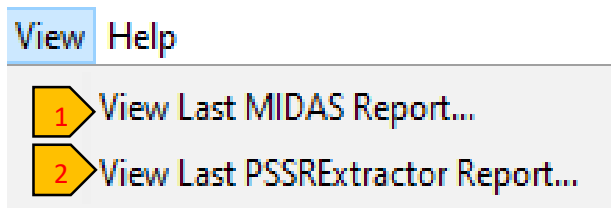
(1) Permite seleccionar el fichero (genoma a escanear en formato FA, FASTA, GBK o GBF ambos de tipo simple o multi-locus) a analizar con MIDAS <sup>1</sup>.

(2) Permite seleccionar el fichero (SSR extraídos con MIDAS en formato MFAA) a analizar con PSSR extractor.

(3) Cierra la aplicación.

(B) Vista:

### Polymorphic SSRs Extractor Pipeline

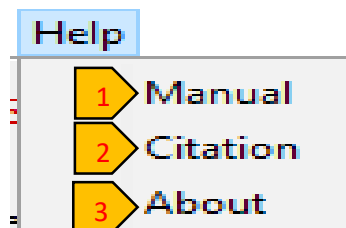


(1) Acceso directo al último reporte de MIDAS.

(2) Acceso directo al último reporte de PSSRExtractor.

(C) Ayuda:

### Polymorphic SSRs Extractor Pipeline



(1) Acceso al manual de usuario.

(2) PDF: Metodología para el minado *in silico* de *loci* polimórficos en microsatélites <sup>2</sup>.

(3) Versión de la aplicación, autores, año de creación y copyright.

MIDAS:

(2) Aplicación para la detección de microsatélites (SSRs) exactos e inexactos.

El TextField (3) permite seleccionar el nombre del fichero (genoma a escanear en formato FASTA o GBFF ambos de tipo simple o multi-locus).

“Max target repeat” (4) fija el tamaño máximo de la unidad repetida a escanear en el genoma.

“Scoring Schema” (5) esquema de parámetros del alineamiento para match, mismatch e indel.

En el TextArea (6) se informará el estado actual de la ejecución de la aplicación (1).

Salidas de MIDAS:

Como salidas MIDAS devuelve tres ficheros de tipo texto que tienen como nombre el fichero de entrada y las extensiones .xls, .dat y .mfaa (el .xls para abrir directamente con Excel u otra aplicación de hojas de cálculo).

Fichero .xls:

NC_021054.gbks - Excel														
Archivo	Inicio	Inserar	Diseño de página	Fórmulas	Datos	Revisar	Vista	LOAD TEST	TEAM	¿Qué desea hacer?	Iniciar sesión	Compartir		
A1	Accession: NC_021054.1													
A	B	C	D	E	F	G	H	I	J	K	L	M	N	
1	Accession: NC_021054.1													
2	Pattern	Length	Copies	Start	End	Score	Matches	Mismatches	Indel	Inaccuracy(%)	5' Flank	5'Entropy	3' Flank	3'Entropy
3	aatagc	6	4	1112818	1112843	52	26	0	0	0	0	1.69	tcaggagacacctcggtttg	1.94
4	accagg	6	3	2414292	2414312	42	21	0	0	0	0	1.68	tcacgatggtgatcgga	1.97
5	acggcc	6	3	3883474	3883491	36	18	0	0	0	0	1.91	caacgccaggcgatcg	1.86
6	acggcg	6	4	1339693	1339718	52	26	0	0	0	0	1.88	ttctgccttgcagagc	1.86
7	atgtcg	6	3	3783028	3783048	42	21	0	0	0	0	1.85	ggtagccctgcagcggg	1.3
8	ccrctg	6	3	3810835	3810853	38	19	0	0	0	0	1.74	tcacggcgccctccaga	1.68
9	acccg	5	3	3112221	3112236	32	16	0	0	0	0	1.46	tccagtcctctctctg	1.72
10	cccg	5	3	3810093	3810107	30	15	0	0	0	0	1.99	acagctcggtcggtgca	1.87
11	cgcg	5	3	1187955	1187973	38	19	0	0	0	0	1.69	tgctctgacatccggcc	1.82
12	acag	4	3	2564288	2564299	24	12	0	0	0	0	1.93	gattgccgagcgtgtgtg	1.87
13	acg	4	3	1930212	1930223	24	12	0	0	0	0	1.94	gcatagagcagggatgga	1.78
14	acg	4	3	3203924	3203935	24	12	0	0	0	0	1.79	cattgcagcgaagatttc	1.99
15	atcc	4	3	2684154	2684166	26	13	0	0	0	0	1.74	ggcttgtgtgtgtgtttga	1.71
16	atgc	4	4	3581621	3581638	27	17	0	1	5.56	cgcaaccggccaagctgatt	1.94	ctggggttcgctgcgggg	1.41
17	atgc	4	3	1775166	1775178	26	13	0	0	0	0	1.86	tcccgatgcagaggtttt	1.95
18	cccg	4	3	2307018	2307029	24	12	0	0	0	0	1.93	tcacgccggatggagcat	1.93
19	cccg	4	3	2497171	2497183	26	13	0	0	0	0	1.97	gcagcgccaacccgatcg	1.74
20	cccg	4	3	2871052	2871065	28	14	0	0	0	0	1.92	gatccagccgcggccag	1.69
21	ccgg	4	3	822580	822591	24	12	0	0	0	0	1.86	tgaggtgttcaggtggcc	1.76
22	ccgg	4	3	1366424	1366435	24	12	0	0	0	0	1.88	ttaccaggttagcgccgc	1.88
23	ccgg	4	3	2729630	2729642	26	13	0	0	0	0	1.95	tgaggtgtgtaacgagtg	1.86
24	ccgg	4	4	2741892	2741907	32	16	0	0	0	0	1.87	cgcaacgcccgacgtgac	1.74
25	ccgg	4	3	3110710	3110721	24	12	0	0	0	0	1.68	cgagtcgggacagctgc	1.91
26	ccgg	4	3	725165	725177	26	13	0	0	0	0	1.88	lugggtgtggtgtgtg	1.77
NC_021054.gbks														

Columnas del Fichero .xls:

**Pattern:** Motivo.

**Length:** Tamaño de la unidad repetida.

**Copies:** Número de copias.

**Start:** Posición inicial en el genoma.

**End:** Posición final en el genoma.

**Score:** Puntuación del alineamiento.

**Matches:** Bases coincidentes.

**Mismatches:** Bases no coincidentes.

**Indel:** Inserciones y supresiones de bases.

**Inaccuracy:** % de inexactitud del repetido, medida de imperfección del mismo.

**5' Flank:** Secuencia flanco al extremo 5'.

**5' Entropy:** Entropía composicional del flanco 5'.

**3' Flank:** Secuencia flanco al extremo 3'.

**3' Entropy:** Entropía composicional del flanco 3'.

Fichero .dat: Presenta en forma no tabular los datos anteriores y permite visualizar el alineamiento de secuencia:

```
Assembly: C:\Users\Alejandro\Desktop\all.gbk\mt
\Mycobacterium_tuberculosis_Beijing_NITR203_uid197218\NC_
021054.gbk
Max Pattern Length= 6
Alignment's Parameters: Match=2 Mismatch=-7 Indel=-7
-----
Accession: NC_021054.1
Found Repeats= 195
-----
Pattern:                aatacg
Exact SSR Pos.:         1112819
Exact SSR Copy Number:  4
Matches:                26
Mismatches:             0
Insertions/Deletions:   0
Inexact SSR Start Pos.: 1112818
Inexact SSR End Pos.:   1112843
Inexact SSR Copy Number: 4
Score:                  52
5' Flank:  gacccggaggccgacccggt  Entropy(0-2): 1.68614
3' Flank:  tcgaggacacctgcggtttg  Entropy(0-2): 1.94188
Aligned SSR:
cgaatacgaatacgaatacgaatacg
|||||
cgaatacgaatacgaatacgaatacg
-----
```

Fichero .mfaa: Presenta los microsátélites detectados en formato multi-fasta, en el cual la región del repetido está marcada en minúscula y los flancos en mayúscula:

El encabezado de este fichero presenta información como el número de acceso del GenBank, el motivo y las posiciones en el genoma.

PSSR Extractor:

(8) Aplicación que permite detectar el polimorfismo en microsatélites (resultados del MIDAS).

El TextField (9) permite seleccionar los resultados del MIDAS (MultiFASTA salida de MIDAS) para analizar, las cuáles serán las consultas(query) en BLAST.

En ComboBox (10) se selecciona el organismo al que pertenecen las secuencias cargadas en 9 el cual a su vez es el organismo en el que BLAST hará la búsqueda de similaridad utilizando las secuencias cargadas en 9 y el genoma del organismo seleccionado en 10.

“Max target sequences”<sup>3</sup> (11) fija el máximo de secuencias similares que BLAST salvará por query (resultados del MIDAS).

El “Expect threshold”<sup>3</sup>(12) filtra secuencias que son menos significativas con un Expect value superior al seleccionado.

(13) El porciento de la palabra en que dos secuencias (de nucleótidos) tienen los mismos residuos en las mismas posiciones en el mismo alineamiento<sup>3</sup>.

“Coverage percent”<sup>2</sup>(14) representa en que porciento dos secuencias (query y subject) tienen la misma cantidad de nucleótidos.

“Use System Proxy” (15), si la conexión del ordenador en el que se está ejecutando este programa tiene proxy entonces esta opción debe ser marcada.

(16) Ejecuta el programa acorde a los parámetros definidos previamente.

(17) Cancela la ejecución del programa. Este botón se activa luego de pulsar “ejecutar” (16).

En el TextArea (18) se informará el estado actual de la ejecución de la aplicación (8).

(19) Salida del programa, este botón cerrará la ventana.

Nota: BLAST posee más parámetros que los expuestos anteriormente, pero para el análisis (análisis de polimorfismo en microsatélites) que se pretende hacer con este programa es necesario que estos se mantengan constantes. A continuación, se exponen dichos parámetros y el valor que tienen por defecto:

**QUERY\_BELIEVE\_DEFLINE:** false.

**DATABASE:** nr.

**LCASE\_MASK:** true.

**FILTER:** F.

**FORMAT\_TYPE:** Tabular.

**PROGRAM:** blastn.

Nota: Las posiciones que brindan start.p y end.p incluyen los flancos que rodean al microsatélite.

Nota: start.p > end.p si subject tiene dirección de 3' a 5' o sea si subject es una secuencia "reverso complemento" puesto que blastn toma en cuenta estas secuencias en su búsqueda de similaridad.

**bp\_num:** Número de pares de bases nitrogenadas entre dos flancos.

**r.units:** Número de unidades repetidas de bases nitrogenadas entre dos flancos.

**direction:** Dirección del subject: "+" si la dirección de la secuencia es de 5' a 3' y "-" si la dirección de la secuencia es de 3' a 5'.

**outliers:** Se representa con "\*" (en caso de que se cumpla, de lo contrario la celda aparece en blanco) los cuales pueden variar de 1 a 5 "\*" y significa que la cantidad de unidades repetidas entre los flancos de los *subjects* es dudosa por ser muy grande, siendo improbable que exista un microsatélite entre ellos.

Nota: Los valores de corte establecidos(límites) para establecer esta excepción fueron mononucleótido: 157 bp, dinucleótido: 364 bp, trinucleótido: 109 bp, tetranucleótido: 45 bp, pentanucleótido: 150 bp y hexanucleótido: 193 bp. Estos valores fueron definidos después de procesar todos los SSR de más de 200 genomas bacterianos, registrando sus tamaños, y estableciendo el corte en 3 veces el rango intercuartil:

\*: Significa que bp\_num > límite.

\*\*: Significa que bp\_num ≥ dos veces el límite.

\*\*\*: Significa que bp\_num ≥ tres veces el límite.

\*\*\*\*: Significa que bp\_num ≥ cuatro veces el límite.

\*\*\*\*\*: Significa que bp\_num ≥ cinco veces el límite.

**exceptions:** muestran etiquetas que corresponden a excepciones:

**"D"** (*degenerated*): Los *subjects* tiene un **Identity Percent** < 90% y/o un **Coverage Percent** < 90%.

**"NF"** (*not found*): No se encontró ningún *subject* en la base de datos con similitud.

**"O"** (*outlier*).

**"U"** (*unpair*): Para una misma secuencia *subject* aparece un flanco y no el otro.

Resultado genérico: Brinda la información relacionada al polimorfismo para cada *query*, es decir para cada SSR.



test10\_generic\_result\_100\_30\_90.0.xls - Excel

Archivo Inicio Insertar Diseño de página Fórmulas Datos Revisar Vista LOAD TEST TEAM ¿Qué desea hacer? Iniciar sesión Compartir

A1 : X ✓ fx #query\_acc\_ver

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
1	#query_acc_ver	access_number	pattern	RUS	RN	inaccuracy	entropy_5	entropy_3	min_RN	max_RN	range	frequency	alleles	PIC	exceptions
2	NC_021054.1 aatagc 1112818-1112843 c:4 s:52 m:25 mm:0 i:0 ina:0 5e:1.69 3e:1.94	NC_021054.1	aatagc	6	4	0	1.69	1.94	3	5	2	0.97	3	0.058	
3	NC_021054.1 accagg 2414289-2414315 c:4 s:44 m:25 mm:2 i:0 ina:7.41 5e:1.77 3e:1.9	NC_021054.1	accagg	6	4	7.41	1.77	1.94	4	5	1	0.01	2	0.02	
4	NC_021054.1 accggc 3883455-3883499 c:7 s:41 m:36 mm:7 i:2 ina:20 5e:1.93 3e:1.69	NC_021054.1	accggc	6	7	20	1.93	1.69	7	8	1	0.005	2	0.01	
5	NC_021054.1 accggg 1339683-1339718 c:6 s:57 m:34 mm:2 i:1 ina:8.11 5e:1.97 3e:1.8	NC_021054.1	accggg	6	6	8.11	1.97	1.86	6	6	0	1	1	0	
6	NC_021054.1 atggcg 3763026-3763048 c:3 s:43 m:21 mm:0 i:0 ina:0 5e:1.85 3e:1.3	NC_021054.1	atggcg	6	3	0	1.85	1.3	3	4	1	0.01	2	0.02	
7	NC_021054.1 cccgga 3810635-3810661 c:5 s:39 m:29 mm:3 i:2 ina:14.7 5e:1.74 3e:1.6	NC_021054.1	cccgga	6	5	14.7	1.74	1.65	5	5	0	1	1	0	
8	NC_021054.1 cccgga 3112218-3112239 c:4 s:34 m:20 mm:2 i:0 ina:9.09 5e:1.46 3e:1.74	NC_021054.1	cccgga	5	4	9.09	1.46	1.74	4	4	0	1	1	0	
9	NC_021054.1 cccgga 3810090-3810123 c:4 s:41 m:28 mm:5 i:0 ina:15.2 5e:1.93 3e:1.92	NC_021054.1	cccgga	5	6	15.2	1.93	1.92	6	7	1	0.01	2	0.02	
10	NC_021054.1 ccggg 1187955-1187976 c:4 s:39 m:21 mm:1 i:0 ina:4.55 5e:1.69 3e:1.86	NC_021054.1	ccggg	5	4	4.55	1.69	1.86	4	4	0	1	1	0	
11	NC_021054.1 accg 2564288-2564299 c:3 s:24 m:12 mm:0 i:0 ina:0 5e:1.93 3e:1.87	NC_021054.1	accg	4	3	0	1.93	1.87	3	3	0	1	1	0	
12	NC_021054.1 accg 1930212-1930223 c:3 s:24 m:12 mm:0 i:0 ina:0 5e:1.94 3e:1.78	NC_021054.1	accg	4	3	0	1.94	1.78	3	3	0	1	1	0	
13	NC_021054.1 accg 3203924-3203935 c:3 s:24 m:12 mm:0 i:0 ina:0 5e:1.79 3e:1.99	NC_021054.1	accg	4	3	0	1.79	1.99	3	3	0	1	1	0	
14	NC_021054.1 atccc 2684154-2684166 c:3 s:26 m:13 mm:0 i:0 ina:0 5e:1.74 3e:1.71	NC_021054.1	atccc	4	3	0	1.74	1.71	3	3	0	1	1	0	
15	NC_021054.1 atgc 3581617-3581638 c:3 s:32 m:20 mm:1 i:1 ina:9.09 5e:1.91 3e:1.41	NC_021054.1	atgc	4	5	9.09	1.91	1.41	5	6	1	0.01	2	0.02	
16	NC_021054.1 atgc 1775166-1775178 c:3 s:26 m:13 mm:0 i:0 ina:0 5e:1.86 3e:1.95	NC_021054.1	atgc	4	3	0	1.86	1.95	3	3	0	1	1	0	
17	NC_021054.1 cccg 2307018-2307029 c:3 s:24 m:12 mm:0 i:0 ina:0 5e:1.93 3e:1.93	NC_021054.1	cccg	4	3	0	1.93	1.93	3	3	0	1	1	0	
18	NC_021054.1 cccg 2497171-2497183 c:3 s:26 m:13 mm:0 i:0 ina:0 5e:1.97 3e:1.74	NC_021054.1	cccg	4	3	0	1.97	1.74	3	3	0	1	1	0	
19	NC_021054.1 cccg 2871052-2871076 c:6 s:30 m:21 mm:4 i:0 ina:16 5e:1.92 3e:1.55	NC_021054.1	cccg	4	6	16	1.92	1.55	6	6	0	1	1	0	
20	NC_021054.1 ccgg 322570-322593 c:5 s:27 m:19 mm:2 i:1 ina:13.6 5e:1.99 3e:1.76	NC_021054.1	ccgg	4	5	13.6	1.99	1.76	5	6	1	0.01	2	0.02	
21	NC_021054.1 ccgg 1366424-1366429 c:3 s:24 m:12 mm:0 i:0 ina:0 5e:1.88 3e:1.88	NC_021054.1	ccgg	4	3	0	1.88	1.88	3	3	0	1	1	0	
22	NC_021054.1 ccgg 2729630-2729642 c:3 s:26 m:13 mm:0 i:0 ina:0 5e:1.95 3e:1.86	NC_021054.1	ccgg	4	3	0	1.95	1.86	3	3	0	1	1	0	
23	NC_021054.1 ccgg 2741892-2741910 c:4 s:33 m:18 mm:1 i:0 ina:5.26 5e:1.87 3e:1.72	NC_021054.1	ccgg	4	4	5.26	1.87	1.72	4	5	1	0.01	2	0.02	
24	NC_021054.1 ccgg 3110706-3110721 c:4 s:35 m:18 mm:0 i:1 ina:6.25 5e:1.68 3e:1.91	NC_021054.1	ccgg	4	4	6.25	1.68	1.91	4	4	0	1	1	0	
25	NC_021054.1 ccgg 725165-725177 c:3 s:26 m:13 mm:0 i:0 ina:0 5e:1.88 3e:1.77	NC_021054.1	ccgg	4	3	0	1.88	1.77	3	3	0	1	1	0	
26	NC_021054.1 ccgg 1492937-1492950 c:3 s:28 m:14 mm:0 i:0 ina:0 5e:1.93 3e:1.86	NC_021054.1	ccgg	4	3	0	1.93	1.86	3	4	1	0.01	2	0.02	
27	NC_021054.1 ccgg 4512374-4512392 c:4 s:31 m:18 mm:0 i:1 ina:5.26 5e:1.94 3e:1.93	NC_021054.1	ccgg	4	4	5.26	1.94	1.93	4	5	1	0.01	2	0.02	
28	NC_021054.1 ccgg 4200878-4200892 c:3 s:30 m:15 mm:0 i:0 ina:0 5e:1.69 3e:1.93	NC_021054.1	ccgg	4	3	0	1.69	1.93	3	4	1	0.01	2	0.02	
29	NC_021054.1 acc 1082253-1082268 c:5 s:27 m:15 mm:1 i:0 ina:6.25 5e:1.91 3e:1.69	NC_021054.1	acc	3	5	6.25	1.91	1.69	5	5	0	1	1	0	
30	NC_021054.1 acc 1677618-1677630 c:4 s:26 m:13 mm:0 i:0 ina:0 5e:1.96 3e:1.85	NC_021054.1	acc	3	4	0	1.96	1.85	4	4	0	1	1	0	
31	NC_021054.1 acc 2814775-2814796 c:7 s:36 m:21 mm:0 i:1 ina:8.7 5e:1.8 3e:1.82	NC_021054.1	acc	3	7	8.7	1.8	1.82	7	7	0	1	1	0	
32	NC_021054.1 acg 4152963-4152978 c:5 s:32 m:16 mm:0 i:0 ina:0 5e:1.97 3e:1.88	NC_021054.1	acg	3	5	0	1.97	1.88	5	5	0	1	1	0	
33	NC_021054.1 acg 4193600-4193614 c:5 s:25 m:14 mm:1 i:0 ina:6.67 5e:1.96 3e:1.88	NC_021054.1	acg	3	5	6.67	1.96	1.88	5	5	0	1	1	0	
34	NC_021054.1 acg 4226629-4226653 c:12 s:34 m:29 mm:8 i:0 ina:21.6 5e:1.94 3e:1.87	NC_021054.1	acg	3	12	21.6	1.94	1.87	12	12	0	1	1	0	
35	NC_021054.1 agc 271321-271333 c:4 s:26 m:13 mm:0 i:0 ina:0 5e:1.44 3e:1.79	NC_021054.1	agc	3	4	0	1.44	1.79	4	4	0	1	1	0	
36	NC_021054.1 acc 1031396-1031411 c:4 s:26 m:13 mm:0 i:0 ina:0 5e:1.85 3e:1.84	NC_021054.1	acc	3	4	0	1.85	1.84	4	4	0	1	1	0	

test10\_generic\_result\_100\_30\_90

Listo

**query\_acc\_ver:** Ver resultado específico.

**access\_number:** Identificador de acceso del query en la base de datos.

**pattern:** Motivo.

**RUS:** Tamaño de la unidad repetida en el query.

**RN:** Número de unidades repetidas en el query.

**inaccuracy:** % de inexactitud del repetido, medida de imperfección del mismo.

**5' entropy:** Entropía composicional del flanco 5´.

**3' entropy:** Entropía composicional del flanco 3´.

**min\_RN:** RN mínimo en el conjunto de query y todos los subjects similares al mismo.

**max\_RN:** RN máximo en el conjunto de query y todos los subjects similares al mismo.

**range:** max\_RN - min\_RN.

**frequency:** Frecuencia alélica que presenta el SSR original(query) a partir del cual se hizo la búsqueda.

**alleles:** Número de alelos encontrados para un *locus* (SSRs con RN diferentes).

**PIC:** Contenido de Información Polimórfica ( $1 - \sum_i p_i^2$ ). Este valor también se conoce en otros contextos como heterocigidad promedio esperada o

diversidad genética de Nei, y da una medida de la probabilidad de que, para un *locus* único, un par de alelos escogidos al azar en la población sean diferentes.

**exceptions:** Muestra todas las etiquetas que corresponden a excepciones (las del resultado específico) en las validaciones del polimorfismo. Hay entradas en el reporte genérico donde pueden aparecer más de una de estas etiquetas pues las excepciones se pueden dar simultáneamente. Si todas las secuencias *subject* presentan excepciones, ya sea de unos o de otras, entonces se colocan las etiquetas de lo contrario la celda aparece en blanco puesto que hay subjects en el reporte específico sin excepciones.