

CEPiNS (Conserved Exon Prediction in Novel Species) - User Manual/Tutorial

This project has a homepage which can be accessed through the following URL:

<http://cepins.org/>

The homepage contains:

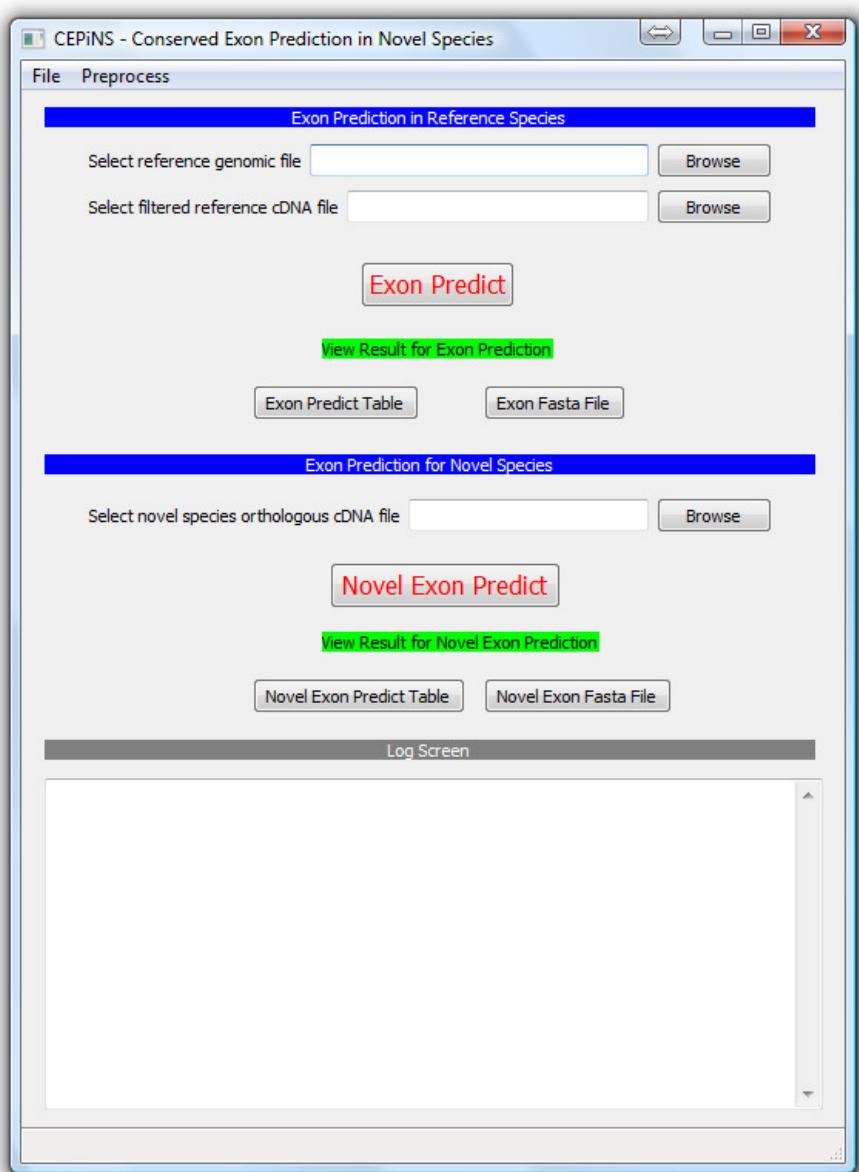
1. Latest version of CEPiNS (Installer which run all versions of Windows)
2. User Manual
3. Example dataset
4. Source code

Additional Information:

This program does not work if your computer has BioEdit installed.

How to use CEPiNS:

After running CEPiNS the user interface is as follows:



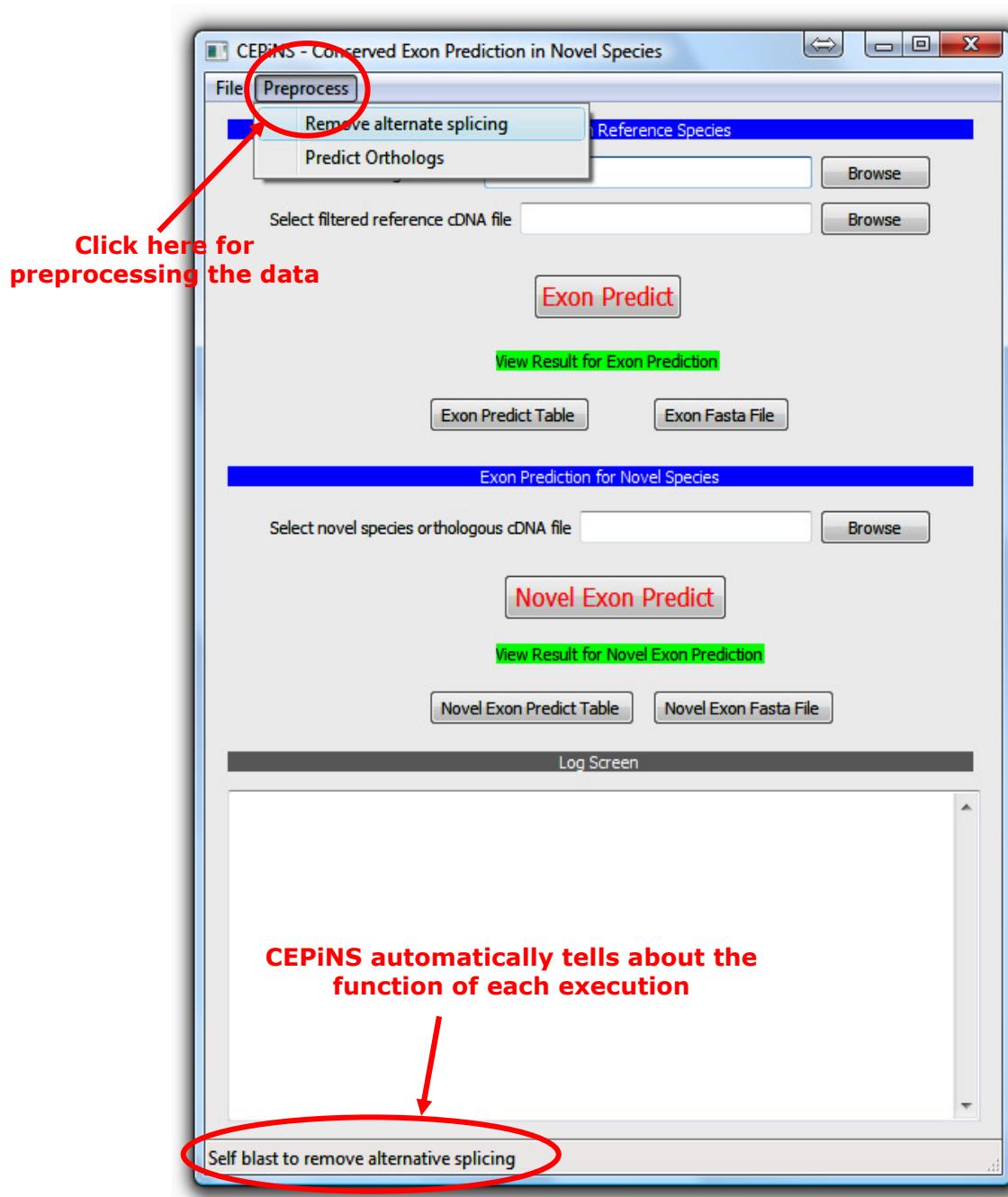
CEPiNS has 3 different parts:

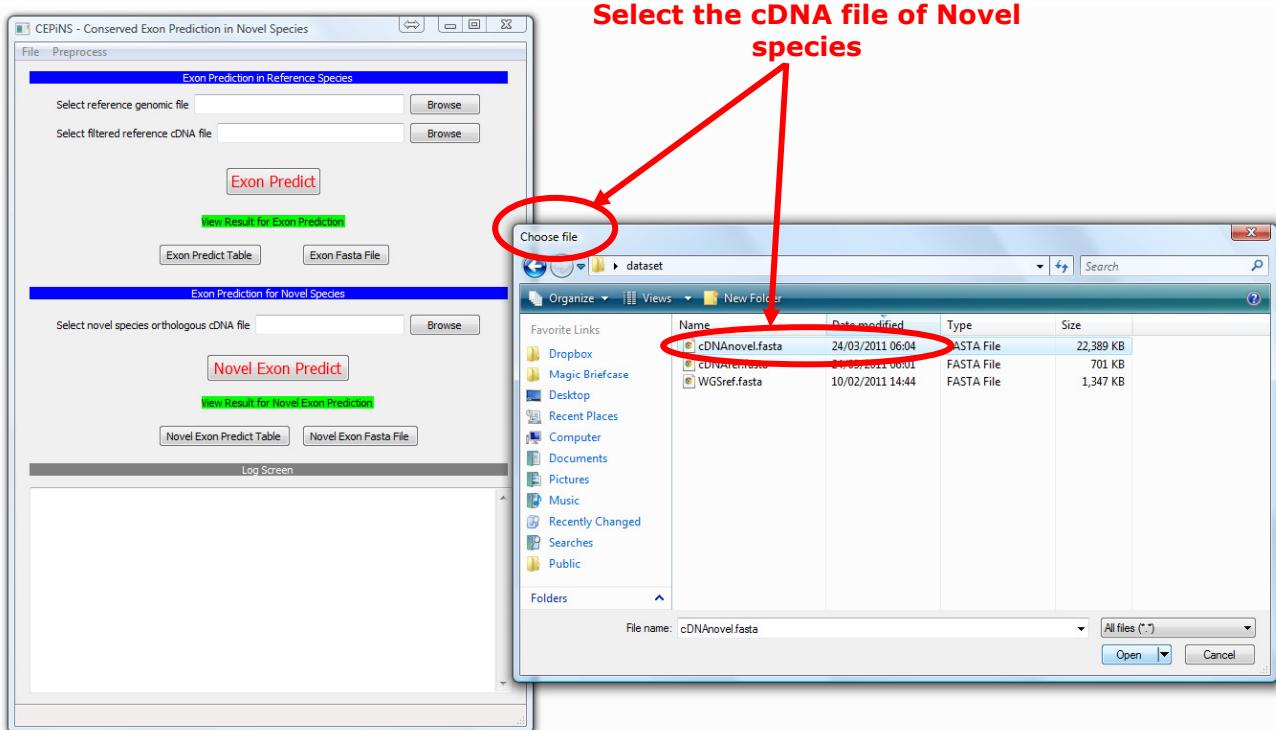
1. Data filtering
 - A. Remove alternate splicing
 - B. Predict Orthologs
2. Exon Prediction in Reference Species
3. Exon Prediction in Novel Species

1. Data filtering

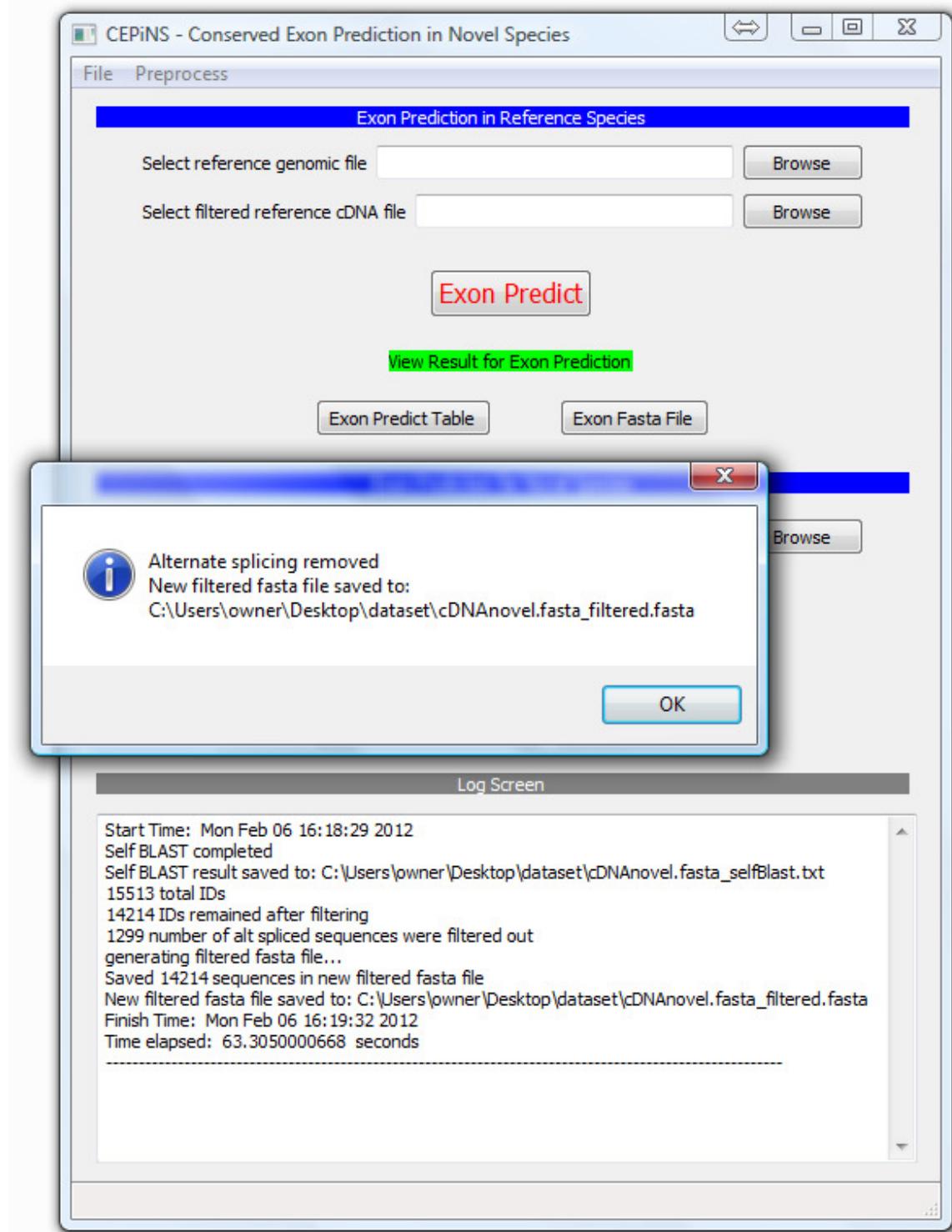
A. Remove alternate splicing:

- I. Click on 'Preprocess' tab from menu bar
- II. Select 'Remove Alternate Splicing'
- III. Select cDNA fasta file
- IV. Click on 'Select'

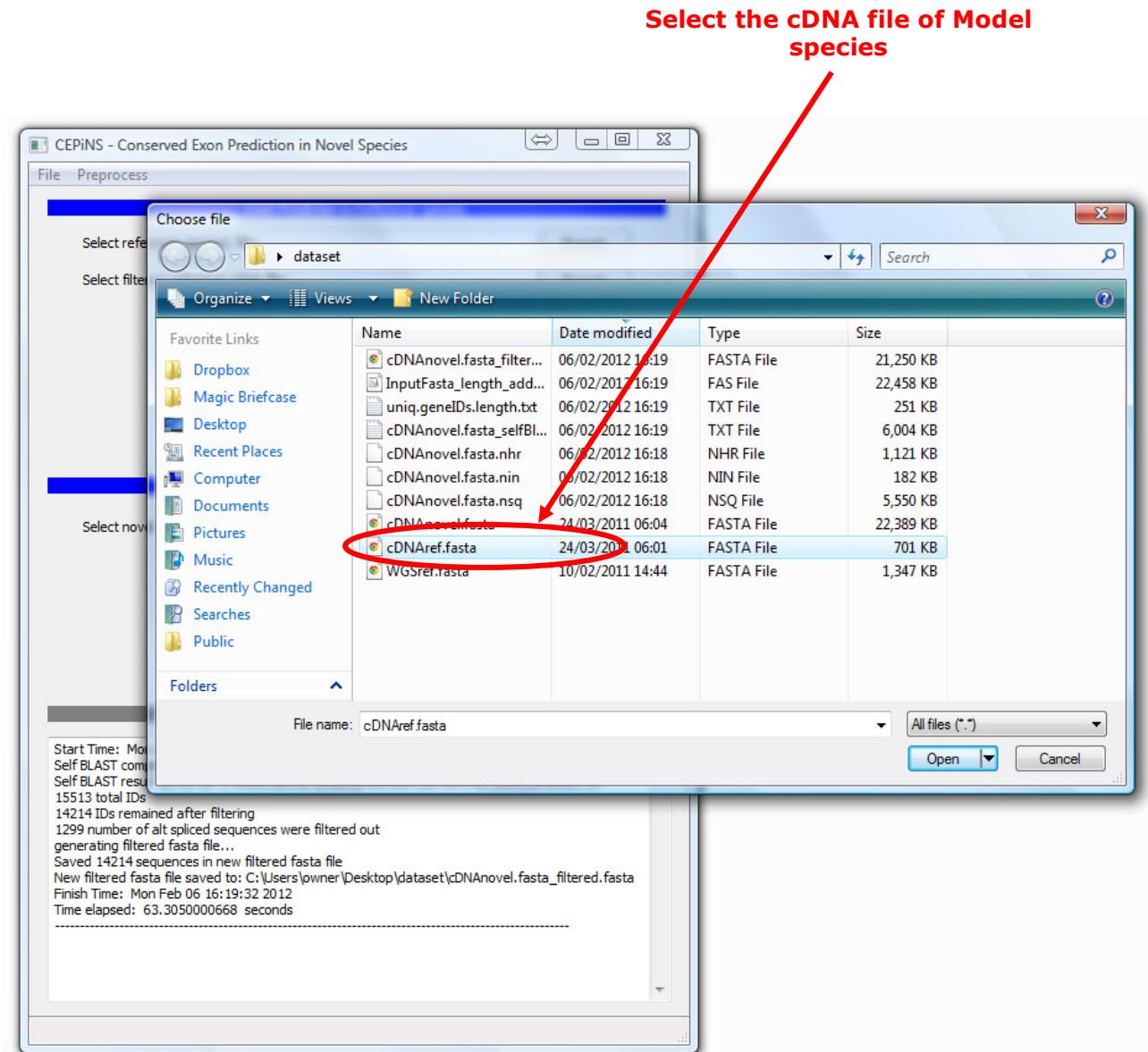


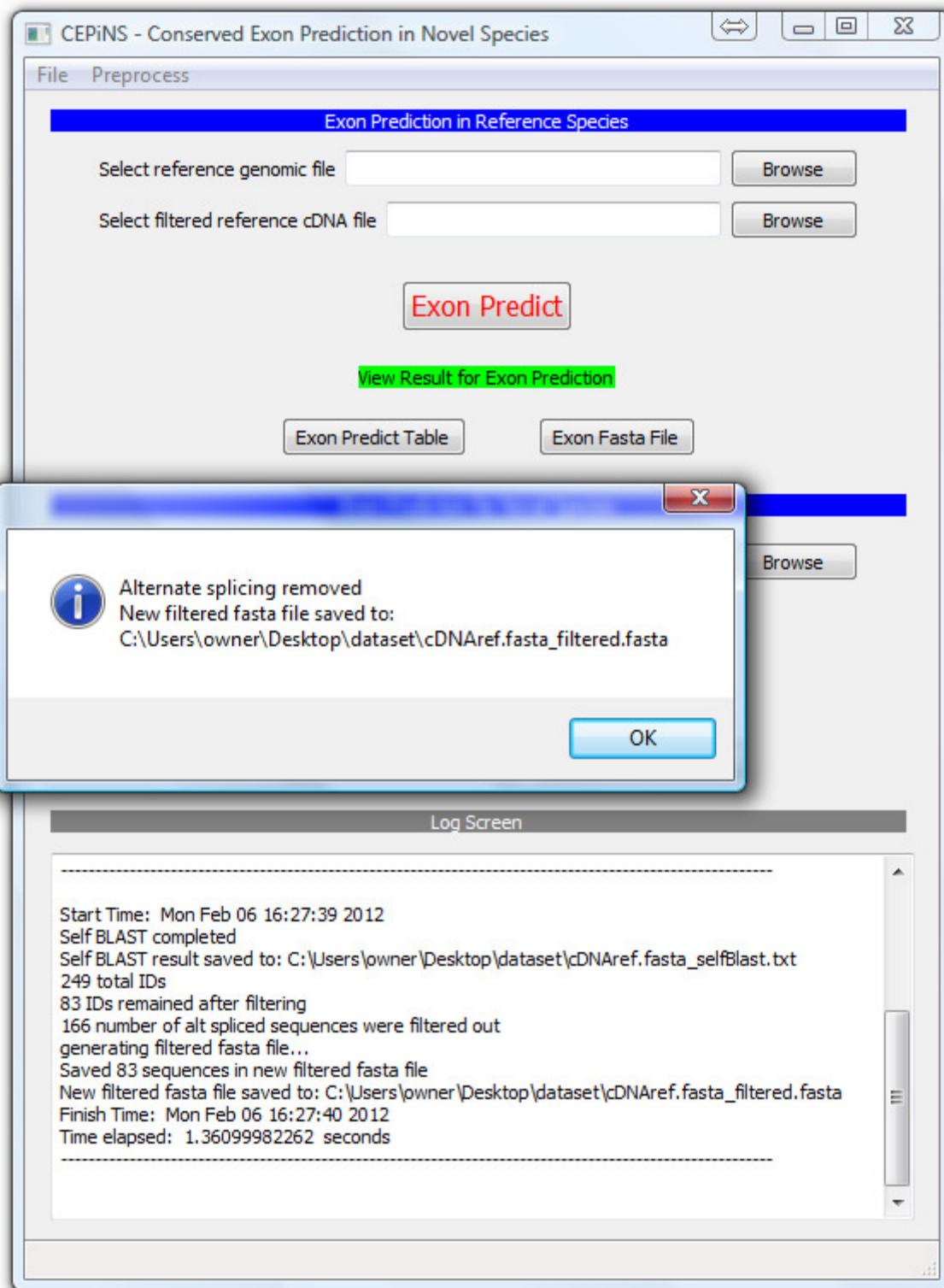


After the execution a new filtered fasta file will be created. The new file name and the saved location will be appeared in a pop up box. The default location is the original dataset location. In the 'Log Screen', all the details will be appeared including the execution time. The 'Log Screen' data can be saved in a text file by clicking 'Save log...' tab from the 'File' menu bar. This log data can be used for keeping records for project's workflow and publication.



Perform the same steps as above for cDNA/gene sequences fasta file of the reference species (model species).

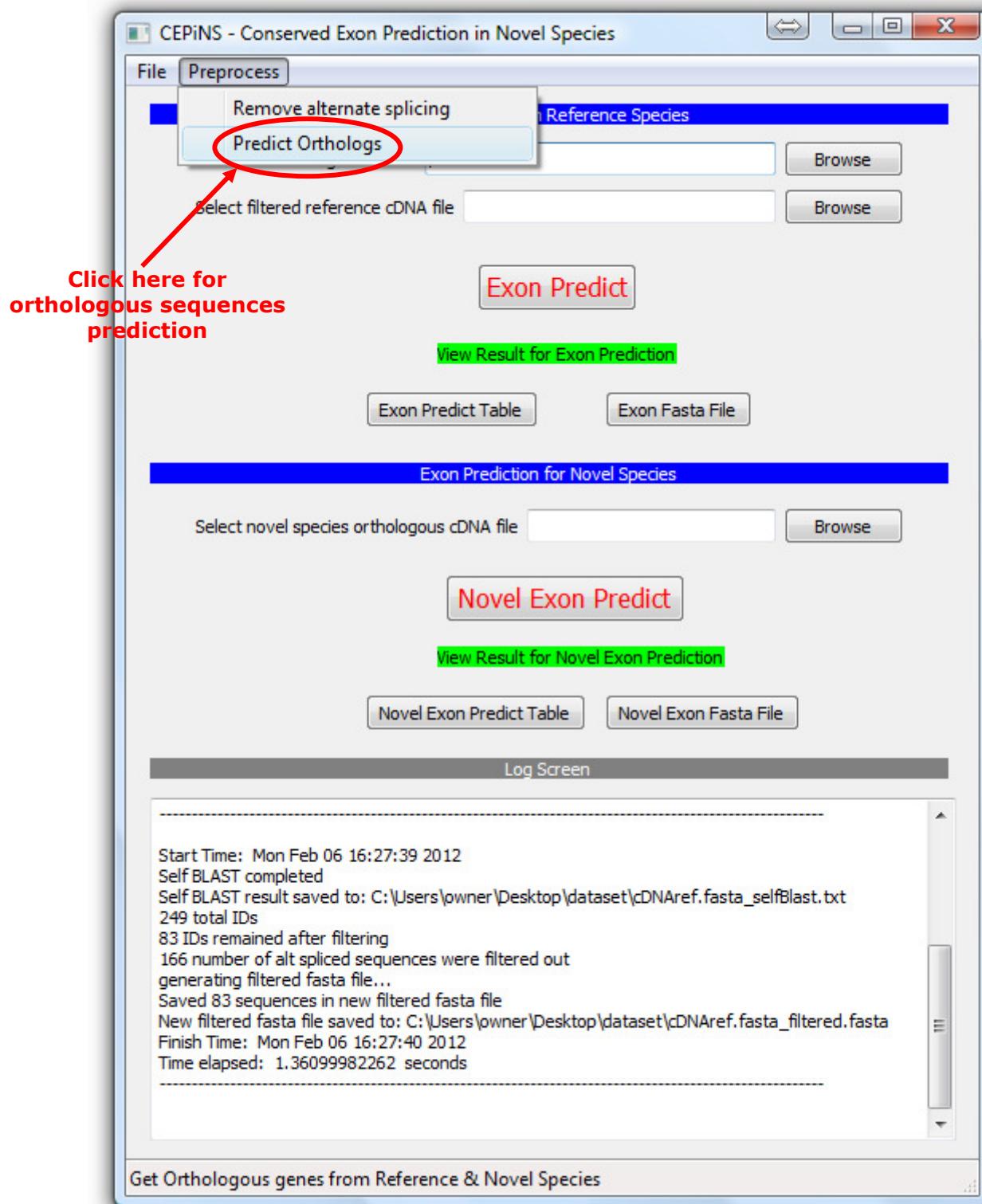




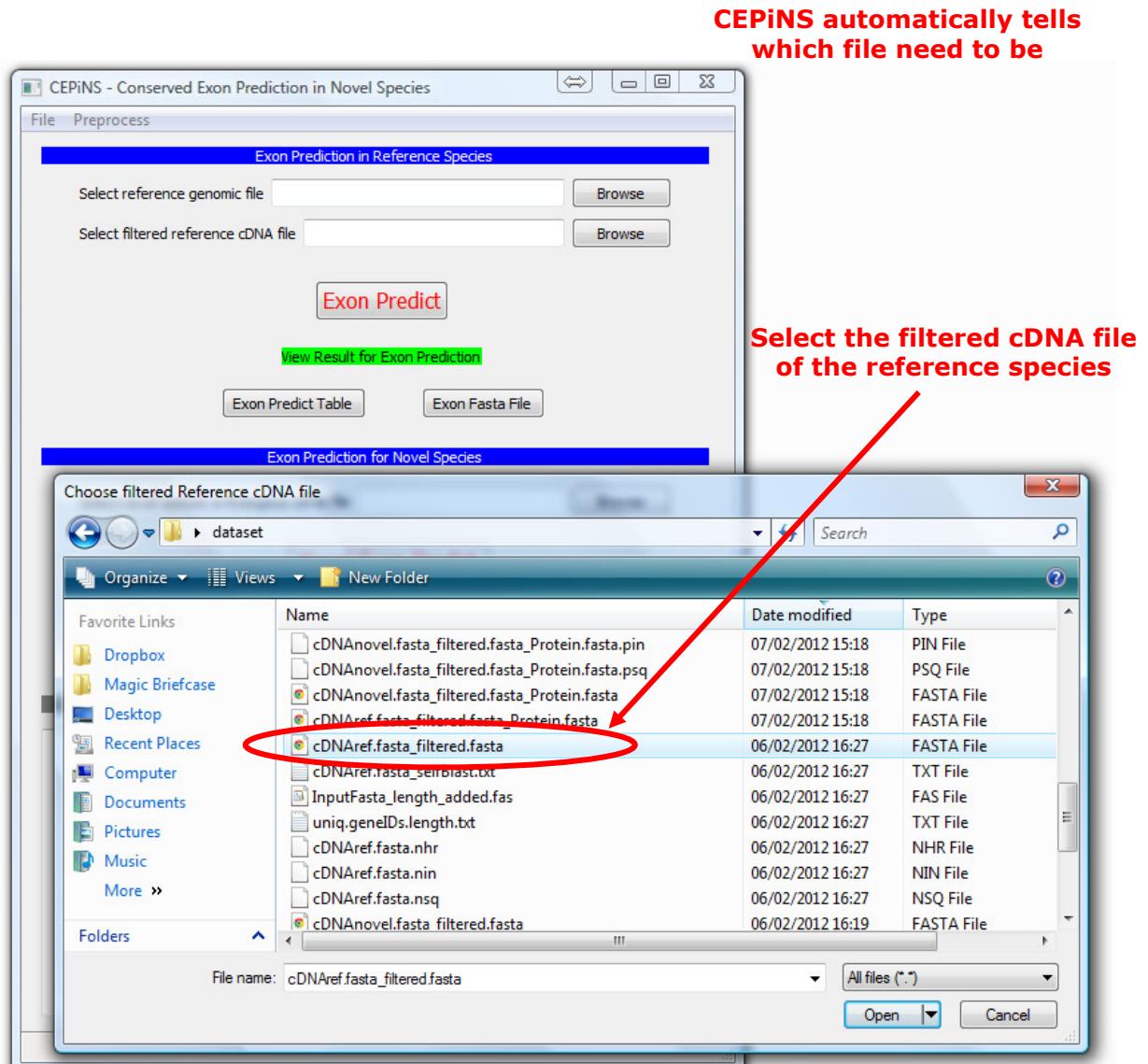
After this step, we will get unique cDNA/gene sequences for both Novel and Reference species where multiple gene copies have been removed. These two files will be used for getting the orthologous sequences.

B. Predict Orthologs:

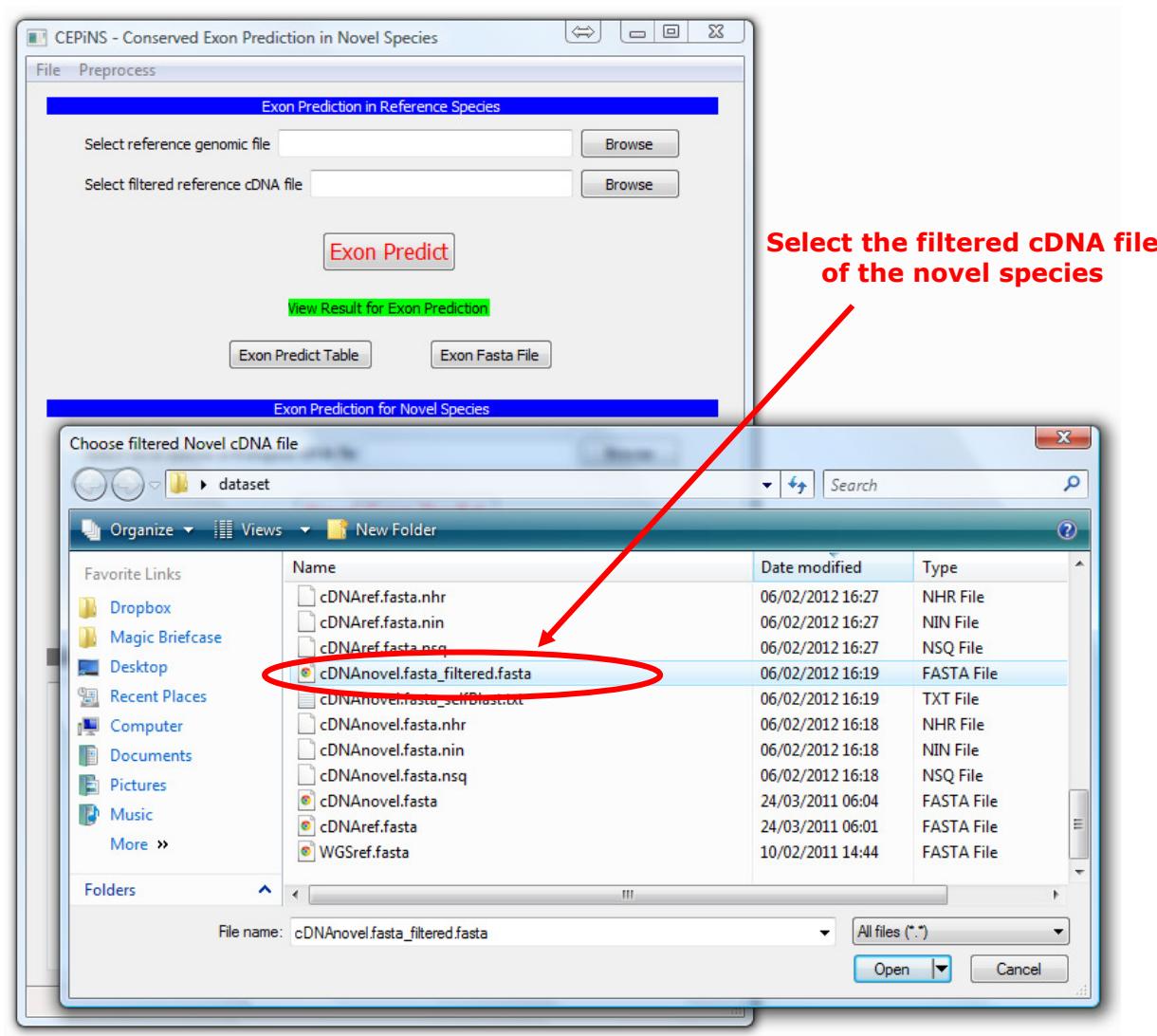
- I. Click on 'Preprocess' tab from menu bar
- II. Select 'Predict Orthologs'
- III. Select filtered novel cDNA fasta file (cDNANovel.fasta_filtered.fasta)
- IV. Click on 'Select'
- V. Select filtered reference cDNA fasta file (cDNARef.fasta_filtered.fasta)
- VI. Click on 'Select'

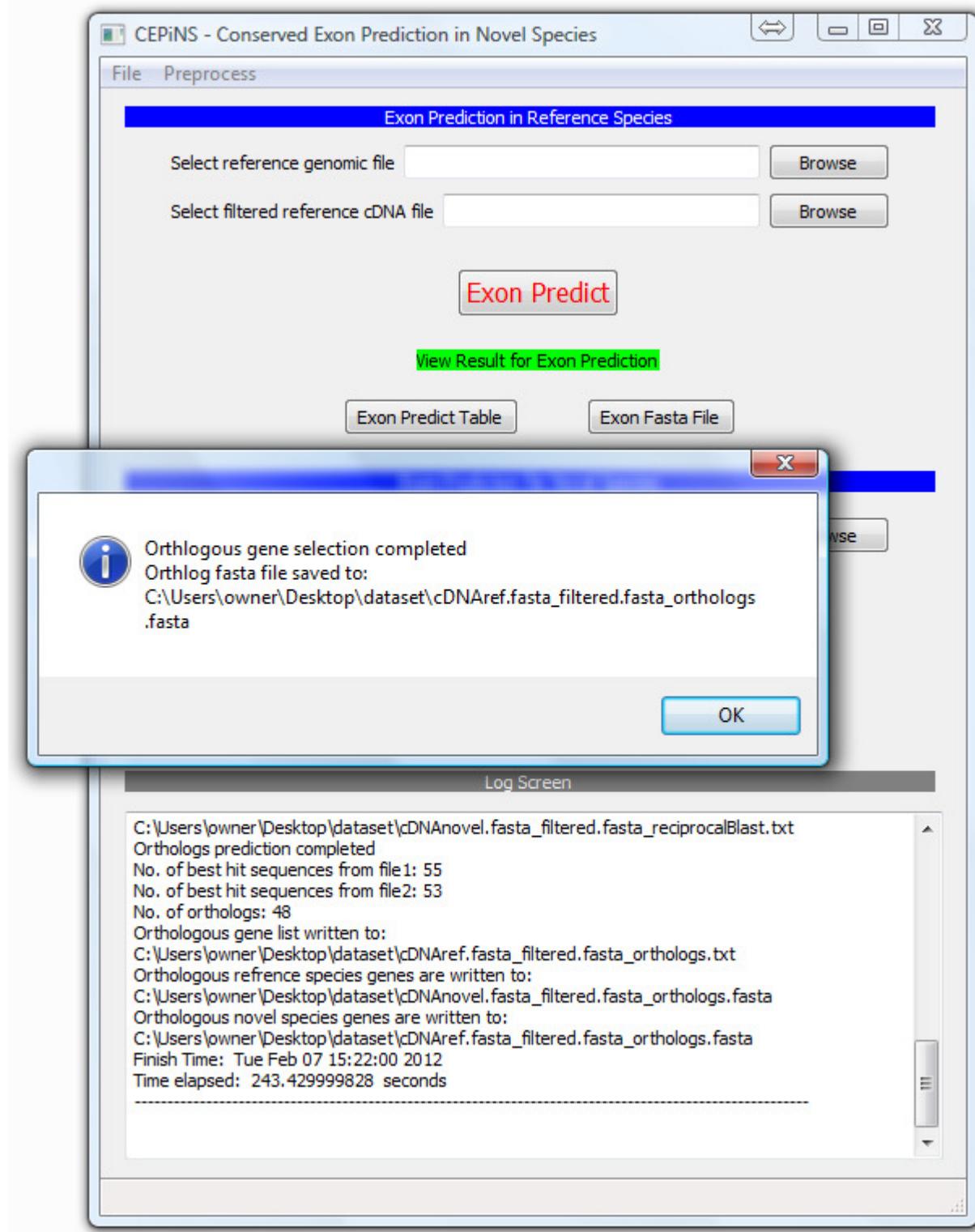


After selecting the 'Predict Orthologs' tab from 'Preprocess' menu bar CEPiNS will be directed to browse option where the filtered cDNA/gene fasta file of reference species is required to be selected.



CEPiNS will be directed to open the filtered cDNA/gene fasta file of novel species.



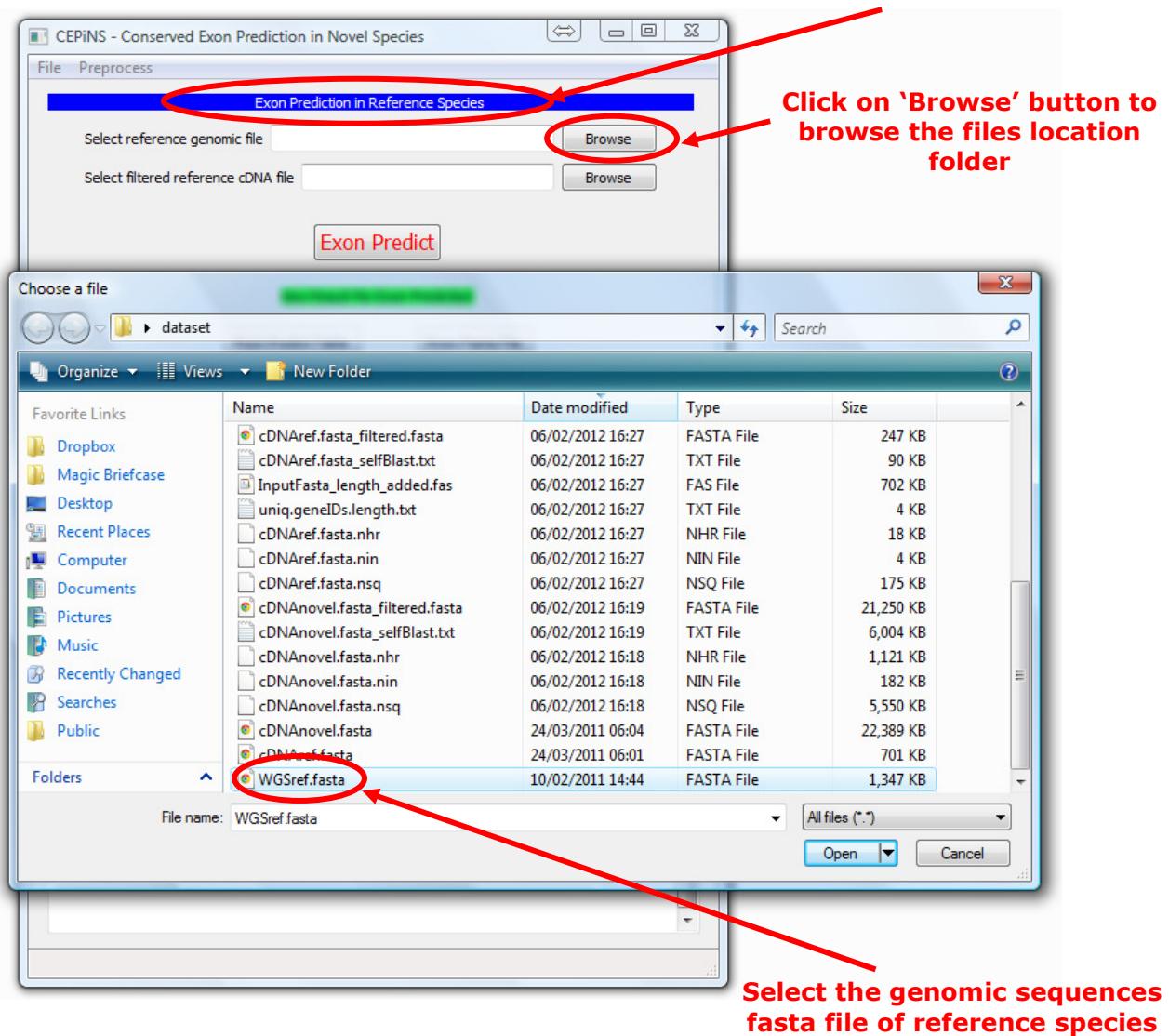


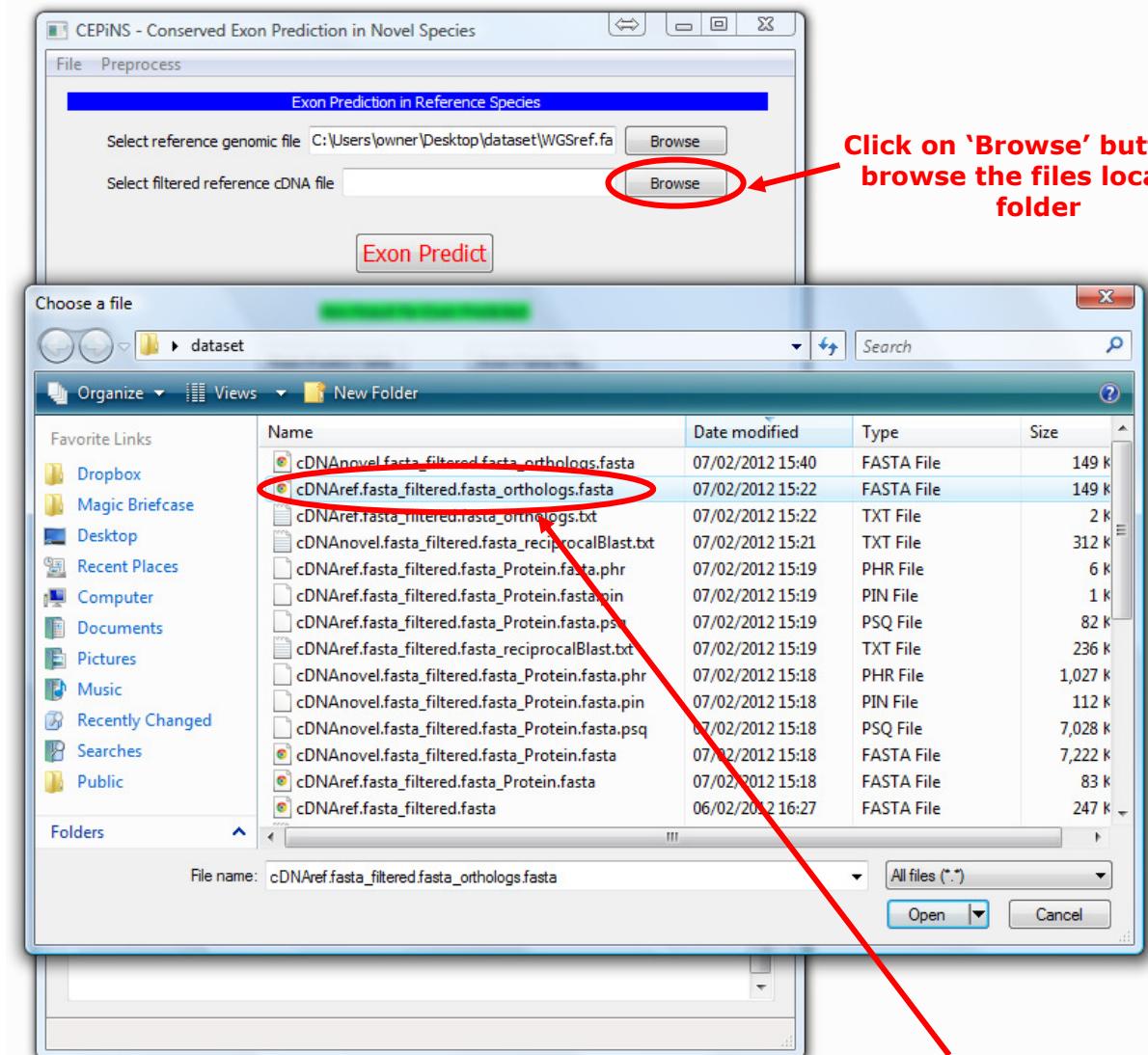
After this step, you will get the orthologous sequences for both reference and novel species. The output files name will be appended according to the original file names where '_orthologs.fasta' will be added at the end of original file names.

2. Exon Prediction in Reference Species:

- I. Select reference genomic file ('WGSref.fasta')
- II. Select filtered reference cDNA file ('cDNARef.fasta_filtered.fasta' or 'cDNARef.fasta_filtered.fasta_orthologs.fasta')
- III. Click on 'Exon Predict' button

Exon Prediction in Reference Species section

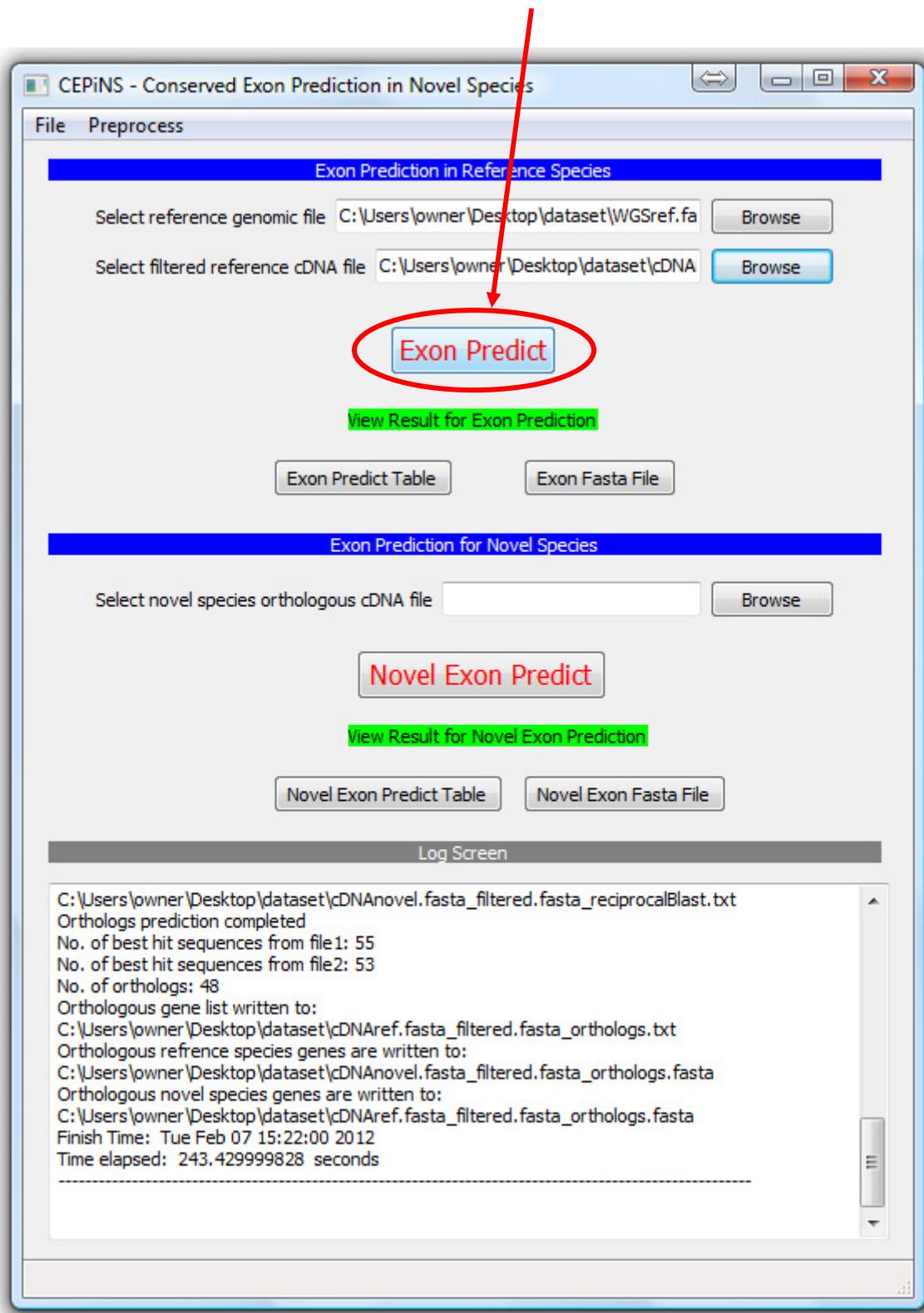


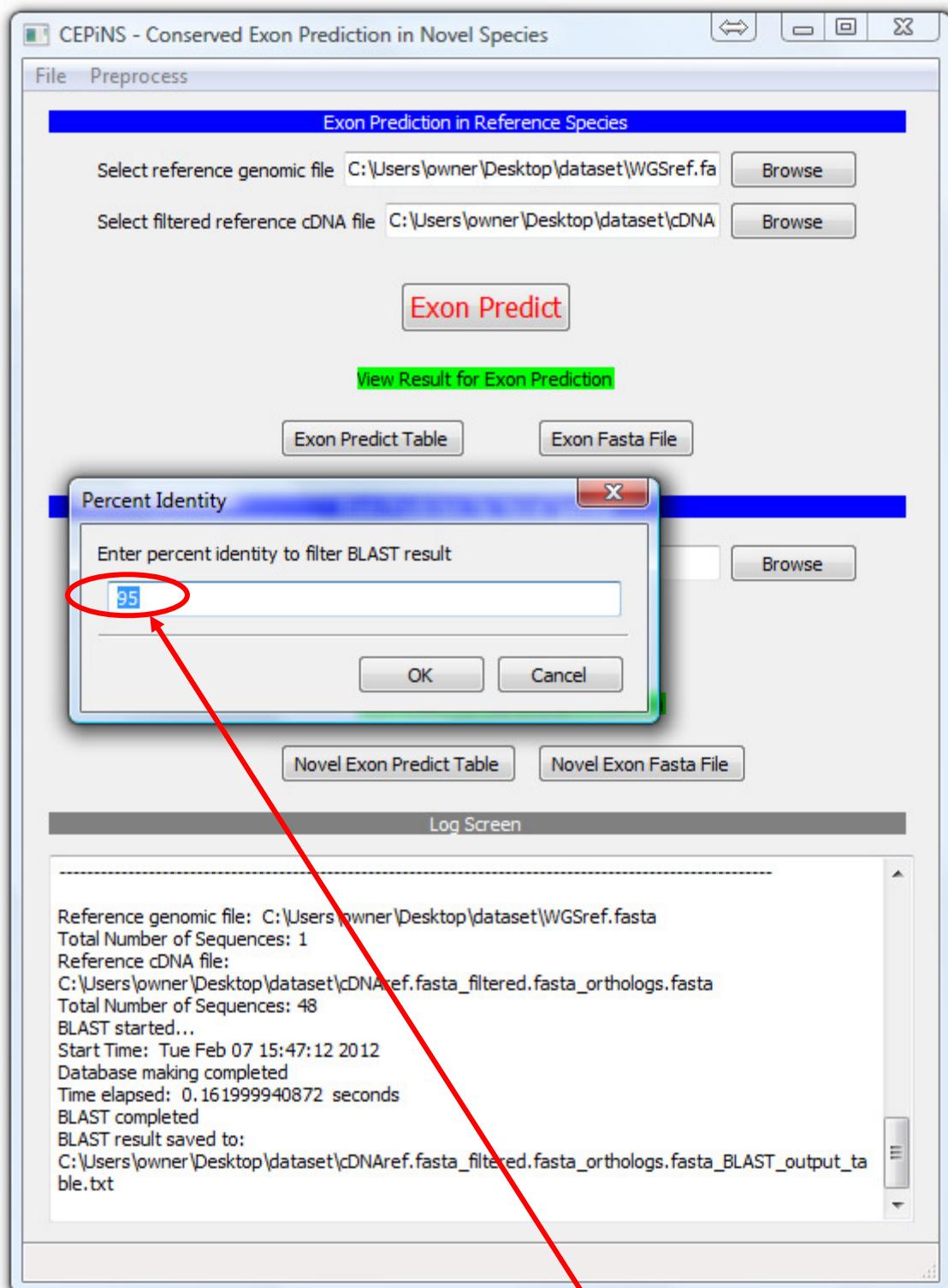


Click on 'Browse' button to browse the files location folder

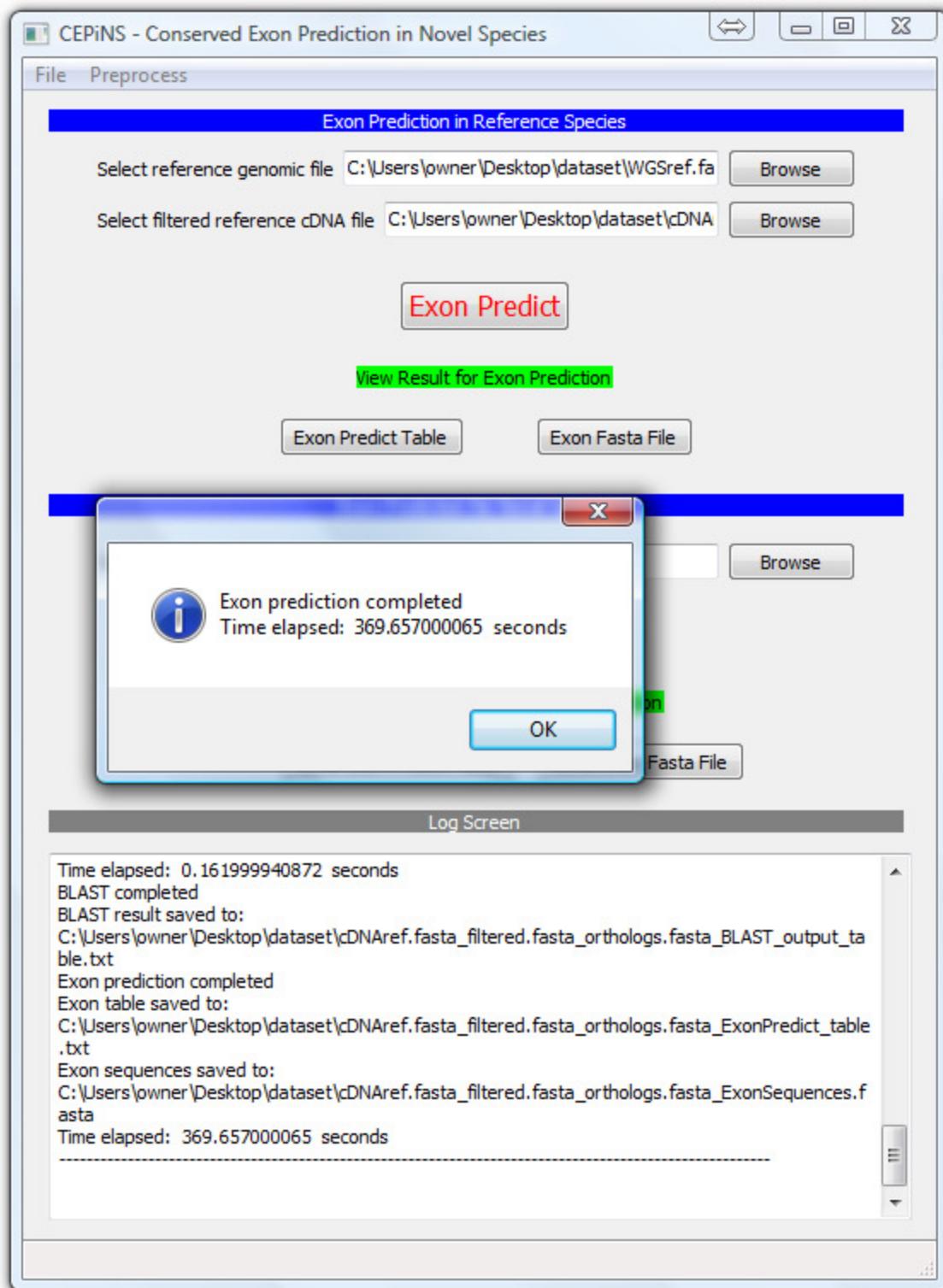
Select the orthologous cDNA/gene fasta file of reference species

Click on 'Exon Predict' button





CEPiNS asks for the percent identity to filter BLAST result which is 95 by default. But user can change the value according to diversity of two species.



After this step, you will get the exon boundaries table text file and a fasta file with the exon sequences which can be opened by clicking 'Exon Predict Table' and 'Exon Fasta File' respectively.

In the exon boundaries table text file, the columns represent cDNA Sequence ID, Location, Exon number, Genomic region location, Location in cDNA, Exon length and Percent Identity respectively.

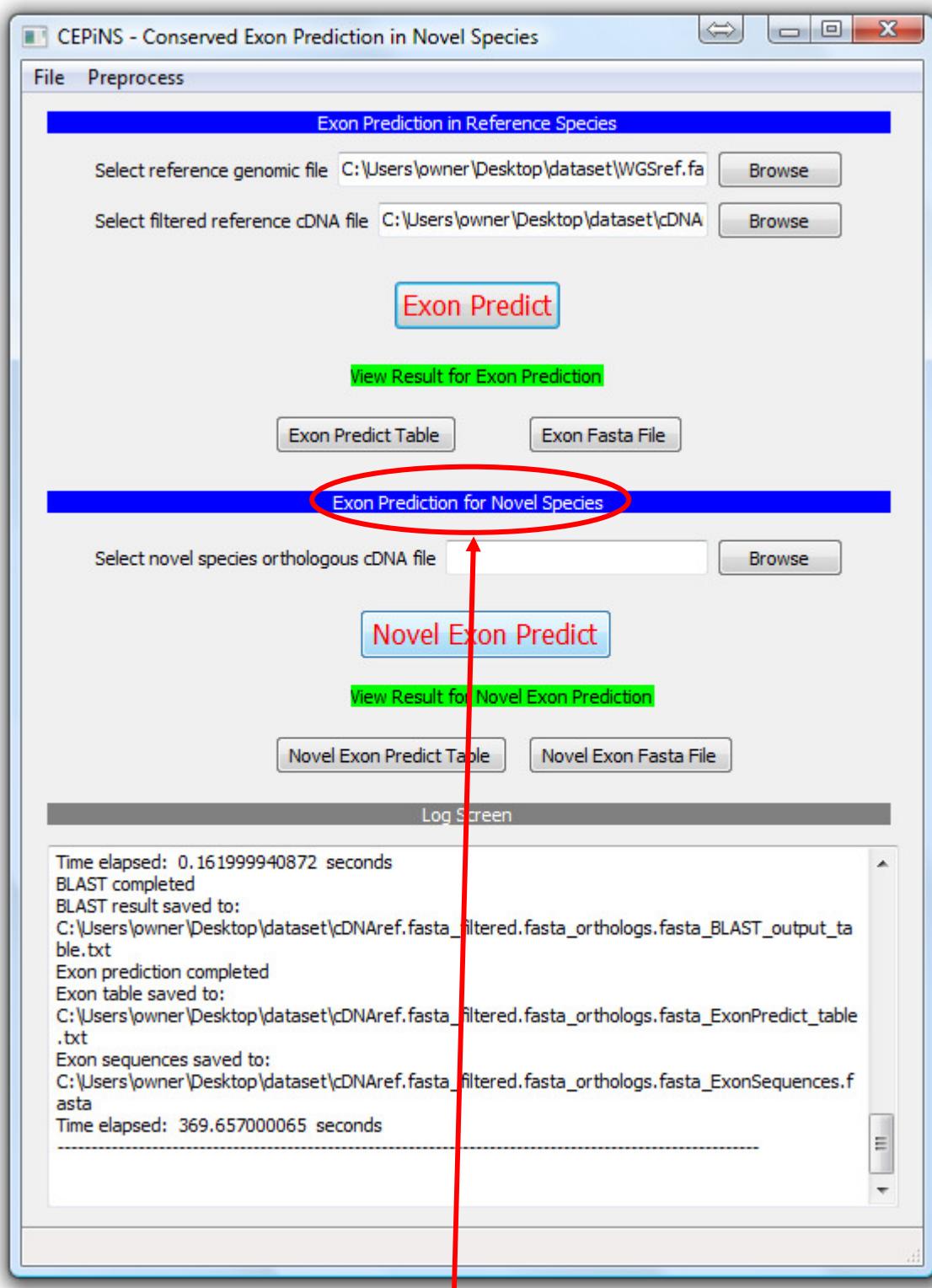
Rf	abg	-PA	*	1	1	1085734-1085809	1-76	76	100.0%
Rf	abg	-PA	*	2	2	1086482-1086511	77-106	30	100.0%
Rf	abg	-PA	*	3	3	1086566-1093731	107-7272	7166	100.0%
Rf	abg	-PA	*	4	4	1093787-1094007	7273-7493	221	100.0%
Rf	abg	-PA	*	5	5	1094065-1096183	7494-9612	2119	100.0%
Rf	abg	-PA	*	6	6	1096242-1096362	9613-9733	121	100.0%
Rf	abg	-PA	*	7	7	1096418-1096740	9734-10056	323	100.0%
Hcf	-PC	*	1	1	1	381317-381632	1-316	316	100.0%
Hcf	-PC	*	2	2	2	382930-383239	317-626	310	100.0%
Hcf	-PC	*	3	3	3	383297-383508	627-838	212	100.0%
Hcf	-PC	*	4	4	4	383572-383760	839-1027	189	100.0%
Hcf	-PC	*	5	5	5	384120-384262	1028-1170	143	100.0%
Hcf	-PC	*	6	6	6	384491-385567	1171-2247	1077	100.0%
Hcf	-PC	*	7	7	7	386300-387389	2248-3337	1090	100.0%
Hcf	-PC	*	8	8	8	390134-390369	3338-3573	236	100.0%
Hcf	-PC	*	9	9	9	390435-390590	3574-3729	156	100.0%
Hcf	-PC	*	10	10	10	391717-391983	3730-3996	267	100.0%
Hcf	-PC	*	11	11	11	392439-392801	3997-4359	363	100.0%
Hcf	-PC	*	12	12	12	392875-393018	4360-4503	144	100.0%
Ephrin	-PA	*	1	1	1	594375-595054	1-680	680	100.0%
Ephrin	-PA	*	2	2	2	595961-596271	681-991	311	100.0%
Ephrin	-PA	*	3	3	3	596328-596950	992-1614	623	100.0%
Ephrin	-PA	*	4	4	4	597604-597948	1615-1959	345	100.0%
sv	-PA	*	1	1	1	1112912-1113254	1-343	343	100.0%
sv	-PA	*	2	2	2	1116202-1116342	344-484	141	100.0%
sv	-PA	*	3	3	3	1116420-1116458	485-523	39	100.0%
sv	-PA	*	4	4	4	1122560-1122725	524-689	166	100.0%
sv	-PA	*	5	5	5	1124161-1124358	690-887	198	100.0%
sv	-PA	*	6	6	6	1126778-1127223	888-1333	446	100.0%
sv	-PA	*	7	7	7	1128325-1128416	1334-1425	92	100.0%
sv	-PA	*	8	8	8	1128724-1128970	1426-1672	247	100.0%
sv	-PA	*	9	9	9	1129100-1129198	1673-1771	99	100.0%
sv	-PA	*	10	10	10	1129735-1130226	1772-2263	492	100.0%
sv	-PA	*	11	11	11	1131407-1131514	2264-2371	108	100.0%
sv	-PA	*	12	12	12	1131914-1132077	2372-2535	164	100.0%

Exon boundaries table of the reference species

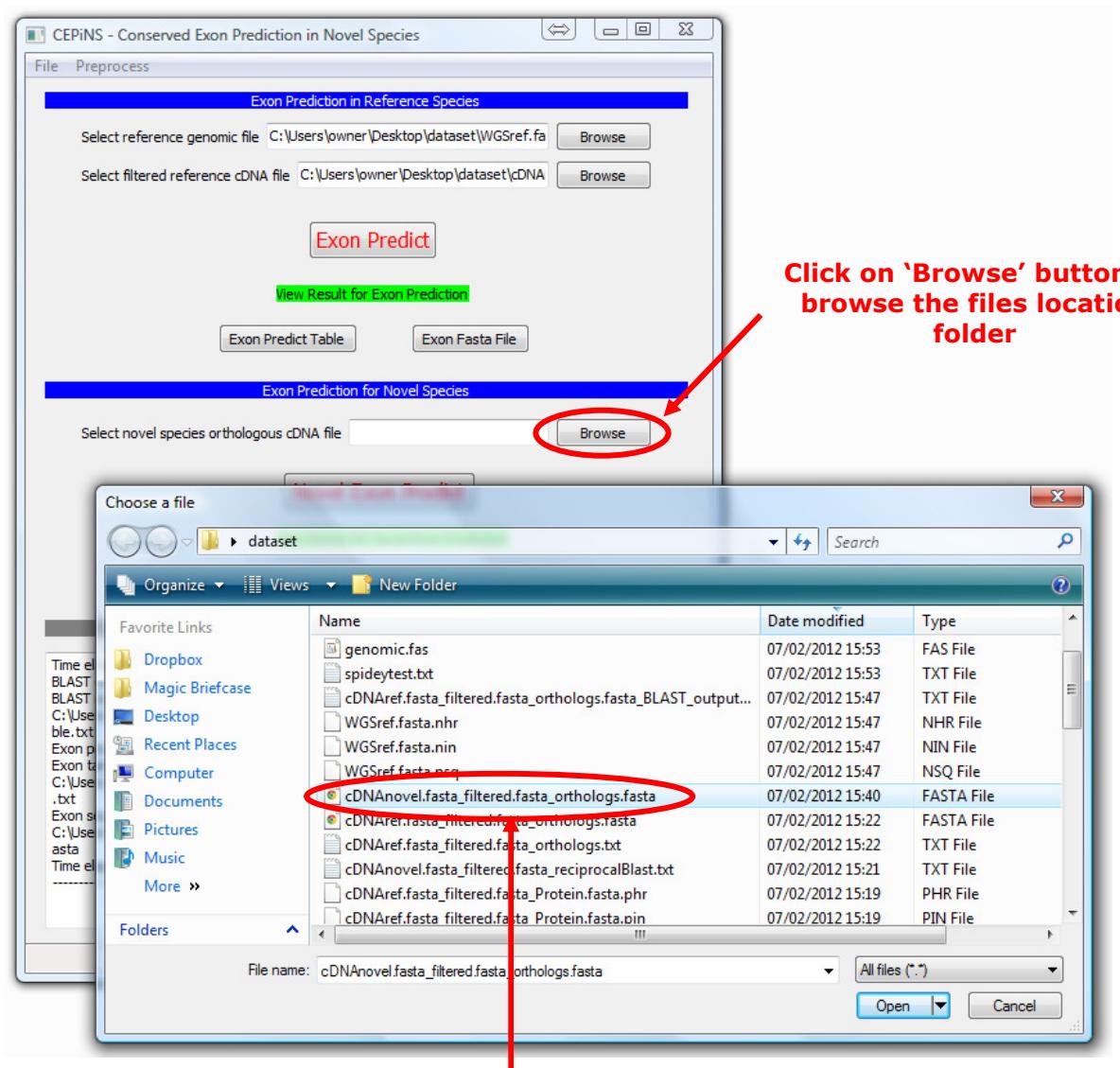
In the exon sequences fasta file, the header contains cDNA Sequence ID, Exon number and Length.

3. Exon Prediction in Novel Species:

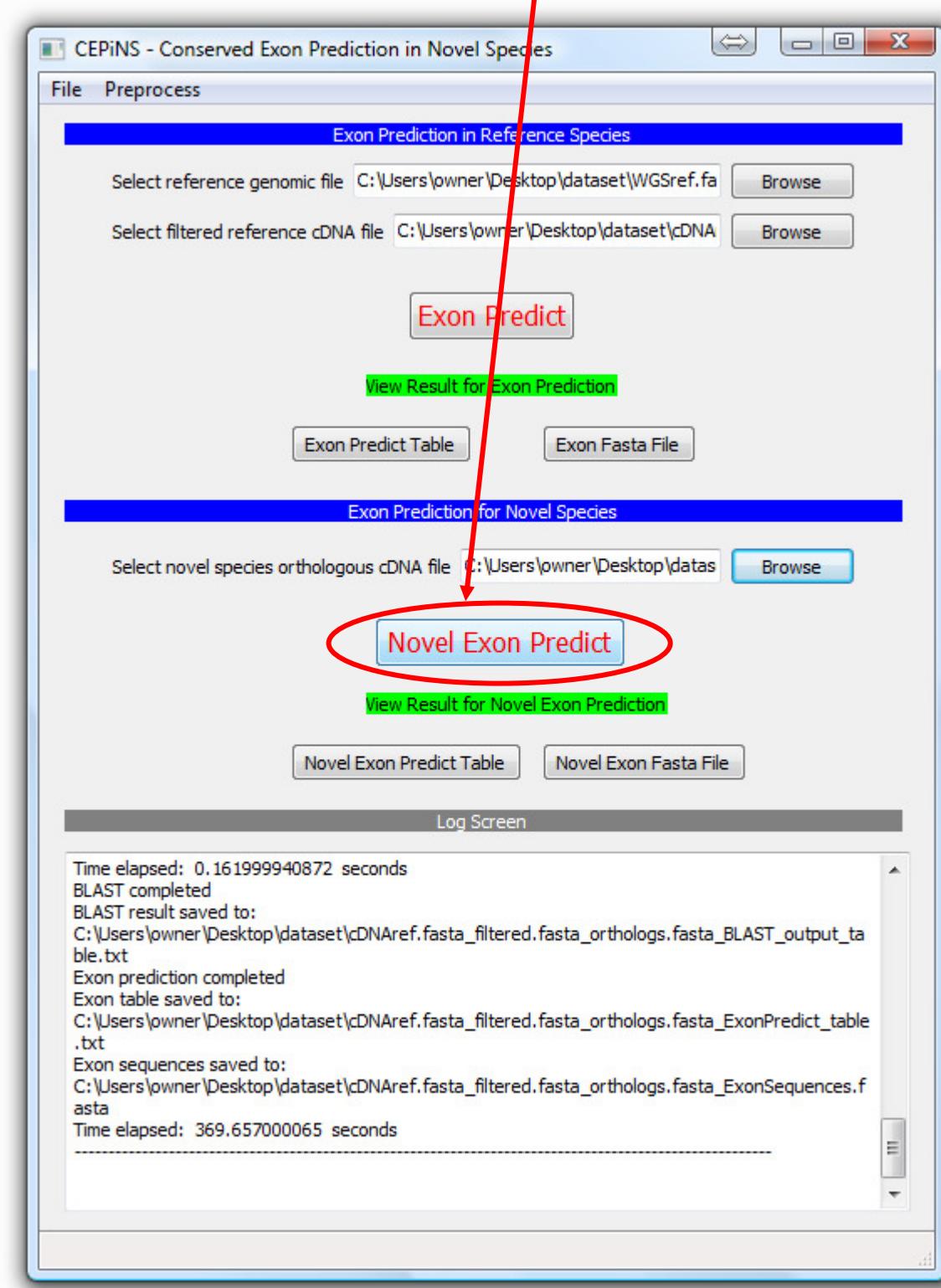
- I. Select novel species orthologous cDNA file
(cDNA novel.fasta_filtered.fasta_orthologs.fasta')
- II. Click on 'Novel Exon Predict' button

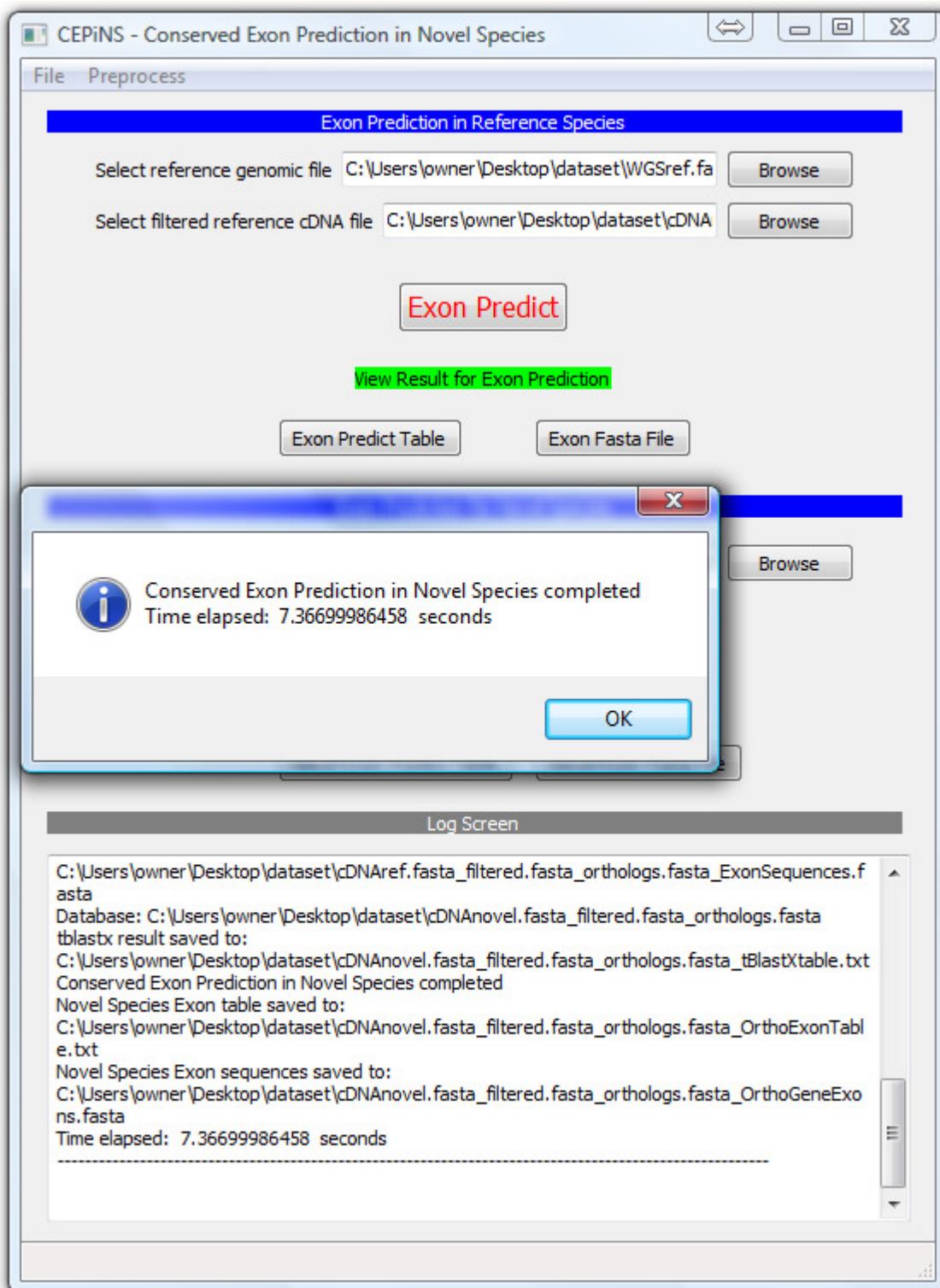


Exon Prediction in Novel Species section



Click on 'Novel Exon Predict' button





CEPiNS writes exon boundaries of the novel species into a tabular text file format and the exon sequences into a fasta file format which can be opened by clicking 'Novel Exon Predict Table' and 'Novel Exon Fasta File' respectively.

After opening the novel exon table text file in Excel, the Excel sheet looks as following.

	A	B	C	D	E	F	G	H	I
1	Reference_ID	Exon_Number	cDNA_coordinate	Length	Novel_ID	Exon_Number	cDNA_coordinate	Length	Identity
2	Rfabg-PA*3	3	107-7272	7166	GK13707-PA	3	107-unk	unk	73.32
3	Rfabg-PA*4	4	7273-7493	221	GK13707-PA	4	7291-7511	221	44.29
4	Rfabg-PA*5	5	7494-9612	2119	GK13707-PA	5	7512-9630	2119	67.99
5	Rfabg-PA*6	6	9613-9733	121	GK13707-PA	6	9631-9751	121	82.5
6	Rfabg-PA*7	7	9734-10056	323	GK13707-PA	7	9752-10074	323	83.18
7	Hcf-PC*1	1	1-316	316	GK13582-PA	1	unk-301	unk	98
8	Hcf-PC*2	2	317-626	310	GK13582-PA	2	302-611	310	98.04
9	Hcf-PC*3	3	627-838	212	GK13582-PA	3	612-823	212	88.57
10	Hcf-PC*4	4	839-1027	189	GK13582-PA	4	824-1012	189	88.71
11	Hcf-PC*5	5	1028-1170	143	GK13582-PA	5	1013-1155	143	91.49
12	Hcf-PC*6	6	1171-2247	1077	GK13582-PA	6	1162-unk	unk	78.26
13	Hcf-PC*7	7	2248-3337	1090	GK13582-PA	7	unk-unk	unk	77.23
14	Hcf-PC*10	10	3730-3996	267	GK13582-PA	10	unk-3975	unk	87.1
15	Hcf-PC*11	11	3997-4359	363	GK13582-PA	11	3976-4338	363	95.87
16	Ephrin-PA*1	1	1-680	680	GK13690-PA	1	unk-842	unk	68.66
17	Ephrin-PA*2	2	681-991	311	GK13690-PA	2	843-1153	311	100
18	Ephrin-PA*3	3	992-1614	623	GK13690-PA	3	unk-unk	unk	76.47
19	Ephrin-PA*4	4	1615-1959	345	GK13690-PA	4	1834-unk	unk	74.19
20	sv-PA*4	4	524-689	166	GK14101-PA	4	41-206	166	81.48
21	sv-PA*5	5	690-887	198	GK14101-PA	5	207-404	198	63.08
22	CG11093-PB*2	2	94-1561	1468	GK13606-PA	2	253-unk	unk	95.19
23	bt-PF*1	1	1-313	313	GK13703-PA	1	1-313	313	93.27
24	bt-PF*3	3	336-498	163	GK13703-PA	3	336-498	163	85.19
25	bt-PF*4	4	499-640	142	GK13703-PA	4	499-640	142	91.49

The exon sequences fasta file of the novel species looks as follows.

```
cDNA novel.fasta_filtered.fasta_orthologs.fasta_OrthoGeneExons.fasta - WordPad
File Edit View Insert Format Help
TGTGAGCAGATTACACCGAAAAGCGACTTTGGACATCTACAAATGCTCTGGAGAGGCTCGTGAAGAG
CGTGGTCCCTGGATCAGATGAGAAGGCCTTCATCTAGCCCTGACTATCCATTCCGCGCTGGCGCAGC
TAAGAGCATCGTTGGTGTTCGAGCGATTCTGAGTACAAGAACCTG
>GK13707-PA_Exon-6 Length:121
TGGAGTTTGCGGGCGCAAATAACAGGAACCTGACCAATTGATGGTGCCTGCTTCATCTGATT
GCTCAGTGAAGGACTGTCTGGAGGGACTGCCAGCCAAAAATTAGTTG
>GK13707-PA_Exon-7 Length:323
GCTTAACTCTCGACTTGTGCCACCACTGAAAACAAGGATAACAAGAACGCCACAAACTCCAGTTG
AGAACGATATGGCATTGACTTTGTTCTAACATGGCGATGGTATTGATACACAGAACTTCGACA
AGTTGAAGGCTCCGGAGCAGAAGAAGGTACTCAACCAAGTTACTTCCTGATTGCCGATACTTGTTCA
AGACCGAAATCGTAGTGAATGCTCTGGCTGGCGCAGCCATGCCCTGATGGTCAGCACAAAGTGTGCCA
TTAAGTCGTCCACCTTCGTACCAAACAAGAACCCAAAGGCCGCGTAG
>GK13582-PA_Exon-1 Length:unk
CGGGATTCGATGAGGCGCTGCTGAACCCAACAGGACCTCAGCCTCGTCTCGTCATGGACATCGCG
CCATCAACATTAAGGAACCTGATGGTTGATTTGGTGGCGTAACGAGGGAAATTGTTGATGAACTTCATG
TCTACAACACAG
>GK13582-PA_Exon-2 Length:310
TTACCAATCAGTGGTATGTGCCGGTGTGAAGGGCGATGTGCCAATGGATGCGCTGCTTACGGATTG
TTGTCGAGGGCACTCGCATTTGCTATTGGCGGAATGATTGAGTATGGAAAAATACTCTAATGAACATT
ATGAGCTCCAGGCCACCAAGTGGAGTGGCGTAAATGTATCCAGAGTCCCCTGACAACGGGTTGTCAC
CCTGICCCCCGCCCTGGGTATAGTTCACCATGGTCGGCGACAAGATACTCCTTTGGTGGCTGGCCA
ATGAATCGGATGATCCGAAAAACATATTCCGAA
>GK13582-PA_Exon-3 Length:212
ATACTGAACGATTGTCACATTTGGACACAGCGGGAGTGCATAGGCCAACAGGCAAGTGGATAATACC
GAAAACGTATGGCGATAGTCGCCGCGCTAGAGAATCGCATACTGGAATCTCGTTACAAGCAAACAC
TGGCAAGCTAAATCTCTGGTTATGGGGATGAGTGGCTGCCGTTGGGAGATTGTTGCTGGGA
For Help, press F1
```