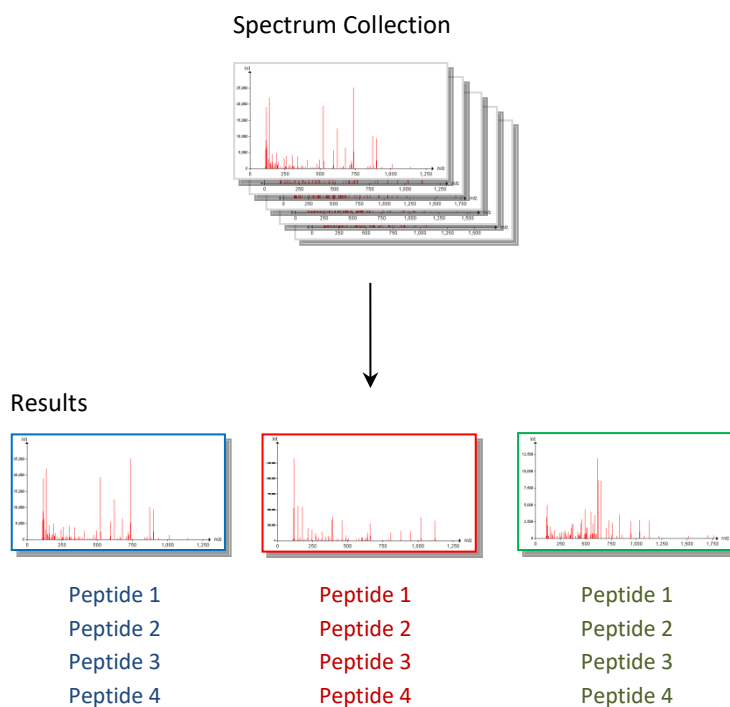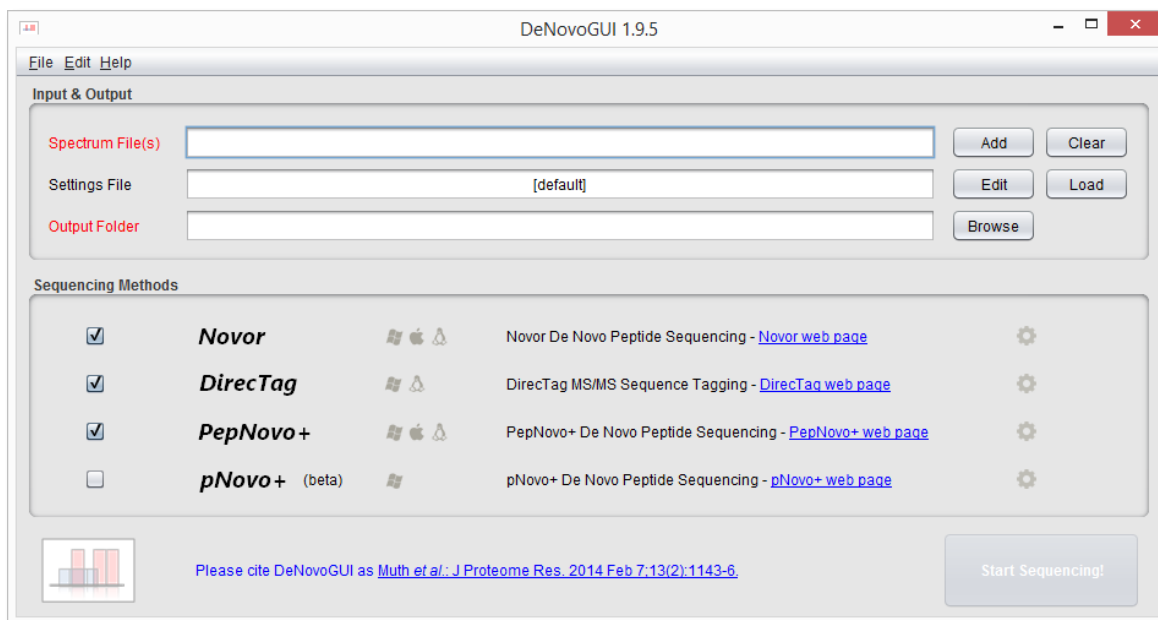# *De Novo* Peptide Identification

When a suitable sequence database is not available for your sample of interest, identification of your acquired fragmentation mass spectra becomes quite difficult. Indeed, the widely-used database search engines cannot be used anymore because they crucially rely on a predefined list of possible protein (and after *in silico* digest, peptide) sequences. In the absence of such a database however, we instead need to rely on the sequence information that can be read directly from a fragmentation spectrum. This procedure is called *de novo* sequencing, and special software tools have been built for this purpose.

We are again going to use the mgf file obtained in the "Peak List Generation" chapter, but instead of searching it against its database (as we did in the "Peptide to Spectrum Matching" chapter), we will attempt to extract as many sequences as possible from the spectra themselves. To do this, we will make use of freely available *de novo* engines. The necessary spectrum file can be found in the resources folder.

Spectrum Collection



Results



| Peptide 1 | Peptide 1 | Peptide 1 |
| Peptide 2 | Peptide 2 | Peptide 2 |
| Peptide 3 | Peptide 3 | Peptide 3 |
| Peptide 4 | Peptide 4 | Peptide 4 |

Novor[1], DirecTag[2], PepNovo+[3] and pNovo+[4] can all easily be used *via* a graphical user interface called DeNovoGUI[5]. DeNovoGUI for Windows platforms is provided in the software folder together with Novor, DirecTag, PepNovo+, and pNovo+. For Mac and Linux versions, please see the DeNovoGUI web page http://compomics.github.io/projects/denovogui.html. Start DeNovoGUI by double clicking the file DeNovoGUI-X.Y.Z.jar (replace X.Y.Z with the current DeNovoGUI version number).

You will then see the following dialog:



The input fields you need to complete are straightforward: provide one or more files with fragmentation spectra, configure some settings, and provide an output folder.

The next section allows you to select the *de novo* tools you wish to apply to your data. You will notice that Novor, DirecTag, and PepNovo+ are already selected. In fact, keen observers may already have noticed these tools in the DeNovoGUI resources folder. This means that when you have downloaded the DeNovoGUI zip file and unzipped it (which comprises the entire installation procedure), you have also already installed all four *de novo* tools along with it!
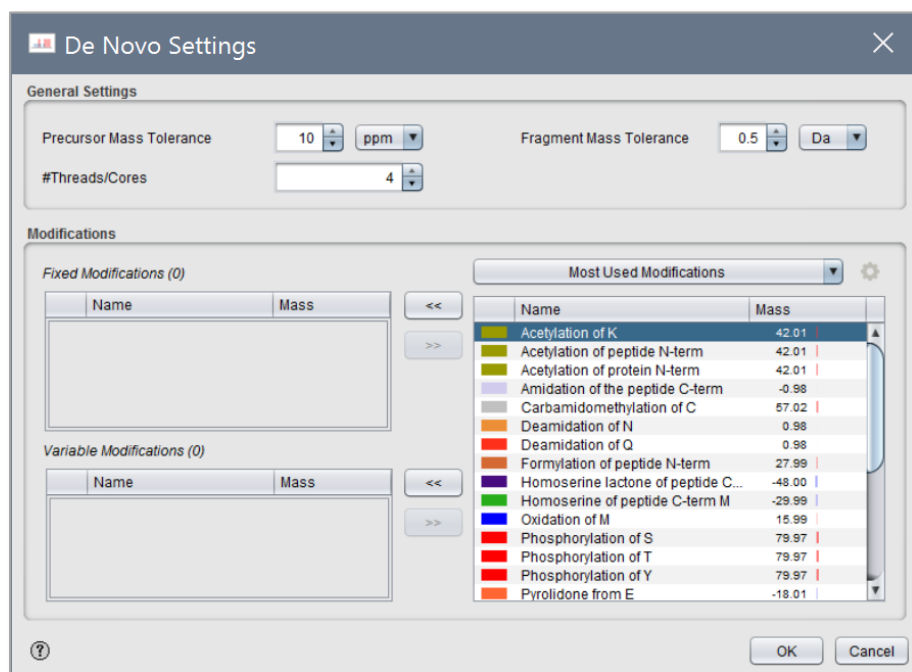
In order to perform the *de novo* identification, we need to provide the spectra and experiment dependent search settings. Load the mgf file qExactive01819.mgf created in the "Peak List Generation" chapter (also available in the resources folder).

---

**Tip:**

*Note that you can load multiple mgf files and even entire folders.*

---

We are now going to set the search settings in the Search Settings dialog. Click the 'Edit' button after the 'Search Settings' text field.



*In which ways is this settings dialog similar to the one presented by SearchGUI (or in fact, any search engine)? In which ways does it differ? Does this make sense?* [1.7a]

Because the data we are using is derived from a high-resolution instrument, we leave the precursor mass tolerance at 10 ppm, and set the fragment mass tolerance to 0.02 Da.

---

*Lennart Martens (lennart.martens@vib-ugent.be) and Harald Barsnes (harald.barsnes@uib.no)*

*Do you think that it is useful in de novo identification to have smaller mass tolerances? Which mass tolerance would be the most important to minimize? [1.7b]*

When using *de novo* it is often possible to set the maximum number of solutions suggested. The default in DeNoboGUI is 10 solutions per spectrum.

*Is it reasonable to expect so many possible sequences? Is this unique to de novo identification, or do you think it affects traditional search engines as well? [1.7c]*

Leave the number of threads at its default; this setting will default to the maximum number of cores on your machine for maximal speed.

The next step is to specify the modifications to consider. As fixed modifications choose Carbamidomethylation of C, and as variable modifications choose Oxidation of M.

> **Tip:**
>
> *'CTRL + Mouse Click' allows you to select multiple entries.*

Note that only the most commonly used modifications are listed in this dialog. There are more modifications available in DeNovoGUI, and you can also set up your own modifications. To see all the modifications, select "All Modifications" in the drop down menu above the modifications table.

To see the modification details or to add your own modifications click the settings ions next to the drop-down menu. (The modification details are also available in the main DeNovoGUI frame, Edit menu > Modifications.) The procedure to manage modifications is exactly the same as in SearchGUI, so we refer the interested reader to the corresponding explanations in the "Peptide to Spectrum Matching" chapter.

All settings are now filled in. Click the 'OK' button and save the settings for future reuse. Choose a meaningful location and name for the parameter file in the dialog. The next time you want to use the exact same search settings you can simply select this file in the main DeNovoGUI dialog.
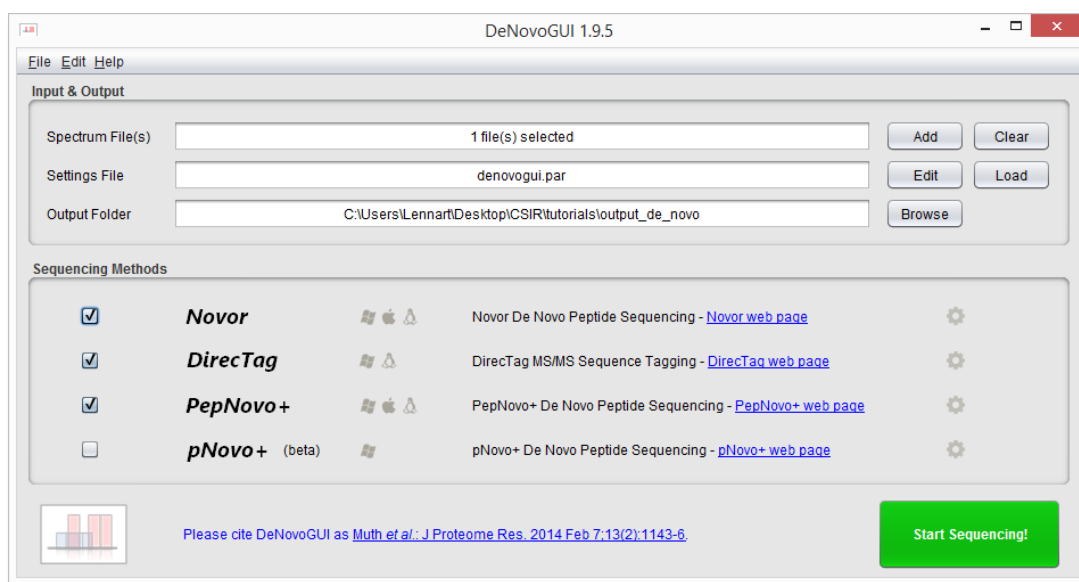
> **Tip:**
>
> *A well-organized library of parameter files can save a lot of time!*

In the main DeNovoGUI dialog you will note that three of the four tools are selected. This means that one can easily run multiple tools and get result files for all of them at the same time. We will keep the default selection.

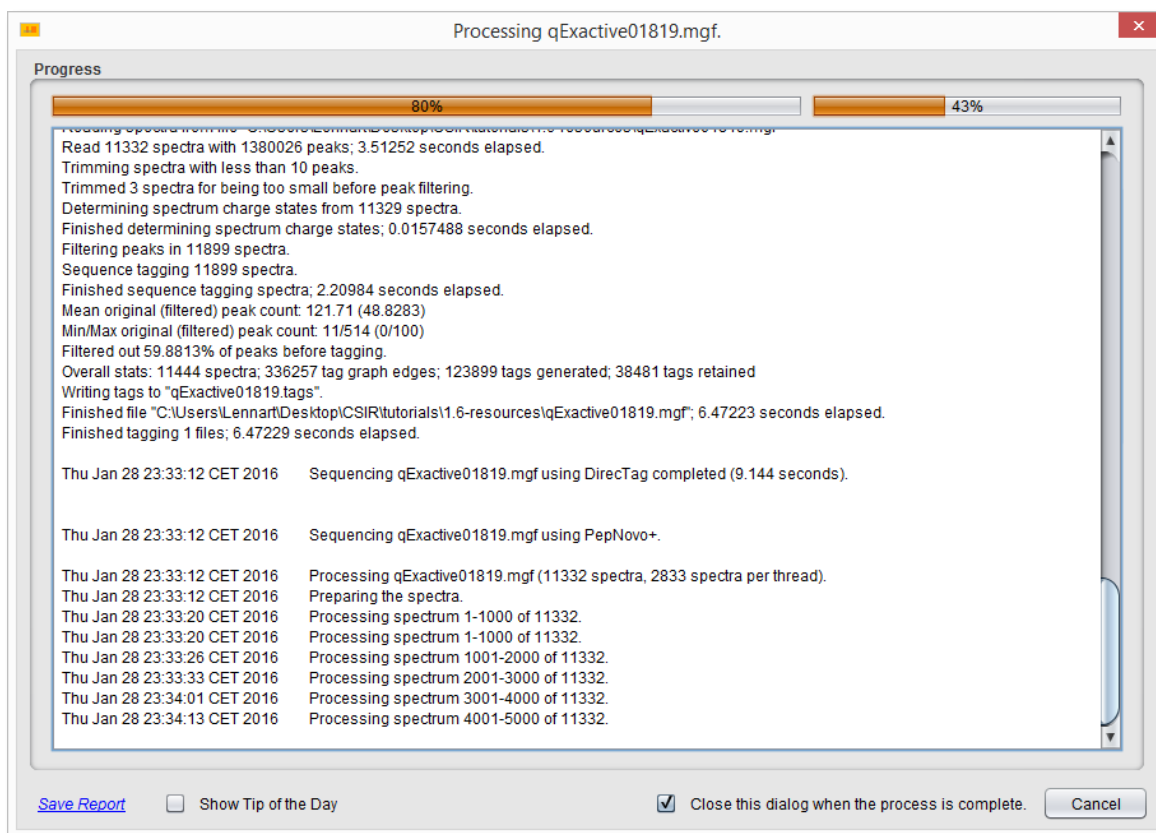Finally, select an output folder. You should now see the following:



---

**Tip:**

*Advanced settings are available for each tool by clicking the settings icon to the right of each tool.*

*Note: changing these advanced settings are recommended for advanced users only!*

---

**Tip:**

*Always use an empty folder for the search output as this simplifies post-processing!*

---

Pressing the 'Start Sequencing!' button will launch the selected *de novo* tools. A progress bar and scrolling text will keep you informed on the progress of the search.

*This work is licensed under the Creative Commons Attribution-ShareAlike 4.0 License.*
*Lennart Martens (lennart.martens@vib-ugent.be) and Harald Barsnes (harald.barsnes@uib.no)*

Page **6** of **11**

Once the *de novo* identification completes, you will see the following results display:



Note that the top panel lists the spectra that were provided as input, along with summary information for each spectrum. Clicking a spectrum then shows the corresponding results in the middle panel. And the bottom panel, finally, shows the annotation of the selected result on the input spectrum.

Now let's inspect this in some more detail. Select the spectrum with the following title "qExactive01819.5053.5053.2 File:"qExactive01819.raw", NativeID:"controllerType=0 controllerNumber=1 scan=5053"". Now have a closer look at the middle, results panel:



There are several rows of results, and these are color coded. At the top, in brown, are the results by DirecTag (revealed when you hover over the colored box). Note that the results are all quite short.

*Why are DirecTag results so short? Is there something special about the way that DirecTag works?* [1.7d]

The orange box coded results are derived from PepNovo+. Note that some of these PepNovo+ results have N-terminal and/or C-terminal gaps, while others do not (notably, the highest scoring and thus first listed peptide sequence is complete, including termini).
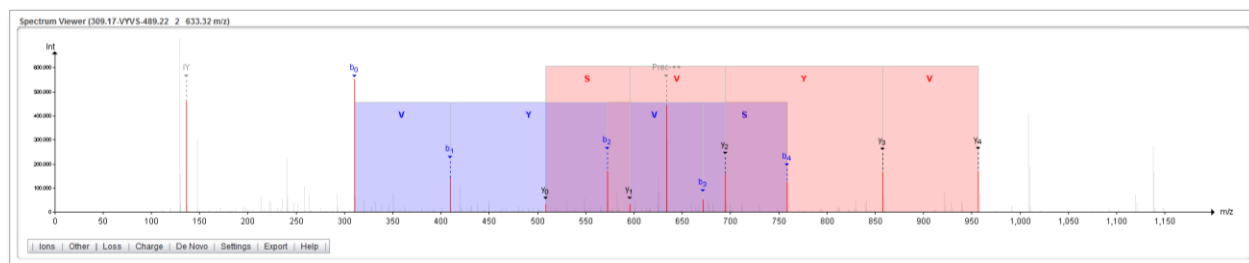
*Where do these terminal gaps come from? In what way does a terminally gapped sequence resemble a sequence tag? [1.7e]*

If you scroll all the way to the bottom, you can see a single blue box coded Novor result.

*Notice the difference between the number of hits provided by Novor, DirecTag and PepNovo+. Does this make sense? Also, do the three tools agree on the potential sequence for this spectrum? [1.7f]*



Now inspect the 6[th] and 10[th] result for this spectrum, which are both sequence tags provided by DirecTag: "309.17-VYVS-489.22" and "507.23-SVYV-291.16", respectively. When you select the 6[th] result, you see the below spectrum matching:
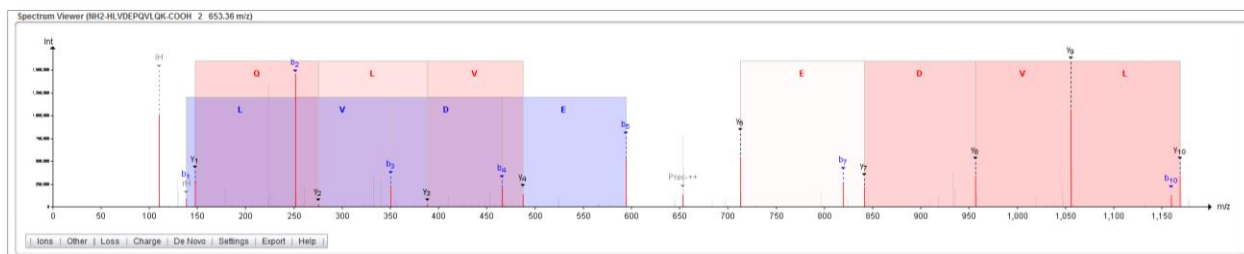


*Why do you see two types of annotation crop up? Which one of these corresponds to the actual tag that we clicked? [1.7g]*

Now select the 10[th] result for this spectrum.

*Compare the two spectrum annotations by selecting both the 6[th] and the 10[th] hit by holding down the Ctrl key. What do you see? Does this make sense? What does this tell you about sequence tags? [1.7h]*

An interesting feature of the Novor algorithm is its residue-specific score. By using the Find feature in the upper right corner of the screen, find the spectrum annotated as HLVDEPQVLQK, and then highlight the blue boxed Novor result in the middle panel for this spectrum. Look closely at the spectrum.



*Do you notice anything special about the annotation? What could be represented?* [1.7i]

There are several ways to work with the results. For individual hits, a BLAST search can be performed quickly and easily by clicking the magnifying glass icon at the far right in the middle table. For instance, clicking this icon for the Novor result that we have just been inspecting, yields this BLAST page:

If you want to perform a more extensive BLAST search, you can also use the 'Export' menu, and then select the 'BLAST' option. This creates a FASTA text file that can be directly submitted to BLAST, in which each peptide is a separate entry. The header for each of these entries contains the *de novo* identification information.

*Why is this link with BLAST useful? Is it a good idea to use a regular BLAST query? Are there any things to keep in mind when doing a BLAST with a de novo sequence?* [1.7j]

There is a further export option that can be very useful, and that is the 'Export', 'Tag Matches' option. This yields a tab-delimited list of all results obtained from your *de novo* identification.

*Bonus exercise for the intrepid reader:*

*Use the 'Tag Matches' export in the tab-delimited format along with a similar export from the original SearchGUI/PeptideShaker identifications from the same spectrum file, and compare the results for each spectrum, and across the whole data set.*

If you encounter any issues with DeNovoGUI, please consult the troubleshooting section at: http://compomics.github.io/projects/denovogui.html.

# References

1.  Ma, B., Novor: real-time peptide de novo sequencing software. *J. Am. Soc. Mass Spectrom.* **26**, 1885-1894 (2015).

2.  Tabb, D.L., Ma, Z., Martin, D.B., Ham, A.L. & Chambers, M.C., Directag: accurate sequence tags from peptide ms/ms through statistical scoring. *J. Proteome Res.* **7**, 3838-3846 (2008).

3.  Frank, A.M., Predicting intensity ranks of peptide fragment ions. *J. Proteome Res.* **8**, 2226-2240 (2009).

4.  Chi, H., Chen, H., He, K., Wu, L., Yang, B. et al., Pnovo+: de novo peptide sequencing using complementary hcd and etd tandem mass spectra. *J. Proteome Res.* **12**, 615-625 (2013).

5.  Muth, T., Weilnböck, L., Rapp, E., Huber, C.G., Martens, L. et al., Denovogui: an open source graphical user interface for de novo sequencing of tandem mass spectra. *J. Proteome Res.* **13**, 1143-1146 (2014).