



# User Manual for pChem

Version 1.0

## Contents

<b>1. Requirements .....</b>	<b>3</b>
<b>2. Download .....</b>	<b>4-5</b>
<b>3. Configuration.....</b>	<b>6-9</b>
<b>4. Run .....</b>	<b>10</b>
<b>5. Output .....</b>	<b>11-15</b>
<b>6. Supporting Protocol 1: Protein sequence database.....</b>	<b>16</b>
<b>7. Supporting Protocol 2: MSconvert.....</b>	<b>17-19</b>
<b>8. Supporting protocol 3: ChemCalc .....</b>	<b>20</b>

# 1. Requirements

## 1) Computing system

pChem search requires a computer with recommended configuration as follows:

- Microsoft Windows 64-bit
- Intel Core i7/i9/Xeon Processor
- 32GB of RAM or more

**Note:** pChem v1.0 is NOT supported by non-Windows operating systems (incl. MacOS, Linux and so on).

## 2) MS Data

- Data dependent acquisition (DDA) with BOTH MS1 and MS/MS spectra recorded in the High-Resolution mode

**Note:** 1) For automatic performance assessment of chemoproteomic probes, it is recommended to acquire MS data from probe-labeled samples with isotope-coding. 2) MS data from non-isotope-labeled samples can also be processed by pChem, but the search result may be subjected to manual inspection.

## 2. Download

- 1) pChem can be freely downloaded from the website:  
<http://pfind.org/software/pChem/index.html>

### pFind Studio: a computational solution for mass spectrometry-based proteomics

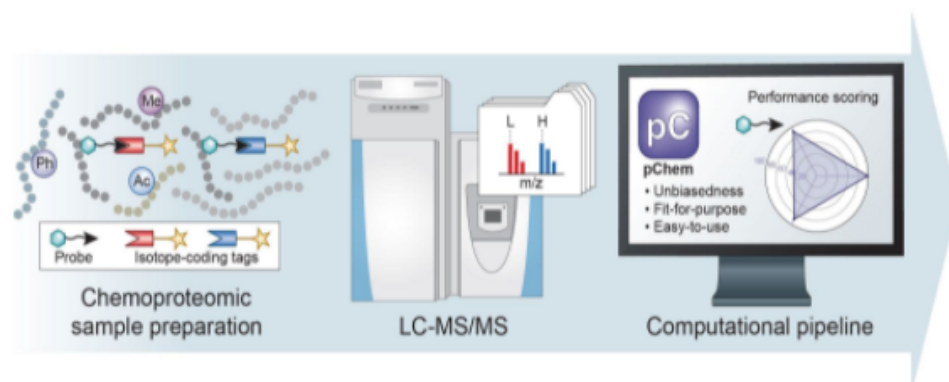
[Home](#) [Members](#) [Publications](#) [Applications](#) [Software](#) [Links](#) [About us](#)

## pChem

[Introduction](#) - [Cite us](#) - [Downloads](#)

### Introduction

Chemical probe coupled with mass spectrometry (MS)-based proteomics, herein termed chemoproteomics, offers versatile tools to globally profile protein features and to systematically interrogate the mode of action of small molecules in a native biological system. Nonetheless, development of an efficient and selective probe for chemoproteomics can still be challenging. Besides, it is also difficult to unbiasedly assess its chemoselectivity at a proteome-wide scale. Here we present pChem, a modification-centric blind search and summarization tool to provide a pipeline for rapid and unbiased assessing of the performance of ABPP and metabolic labeling probes. This pipeline starts experimentally by isotopic coding of PDMs, which can be automatically recognized, paired, and accurately reported by pChem, further allowing users to score the profiling efficiency, modification-homogeneity and proteome-wide residue selectivity of a chemoproteomic probe.



### Cite us

pChem: a modification-centric assessment tool for performance of chemoproteomic probes.  
Ji-Xiang He, Zheng-Cong Fei, Fu-Chu He, Si-Min He, Hao Chi, Jing Yang.  
Under review

### Downloads

pChem version 1.0 is currently free to use. [click to download](#).

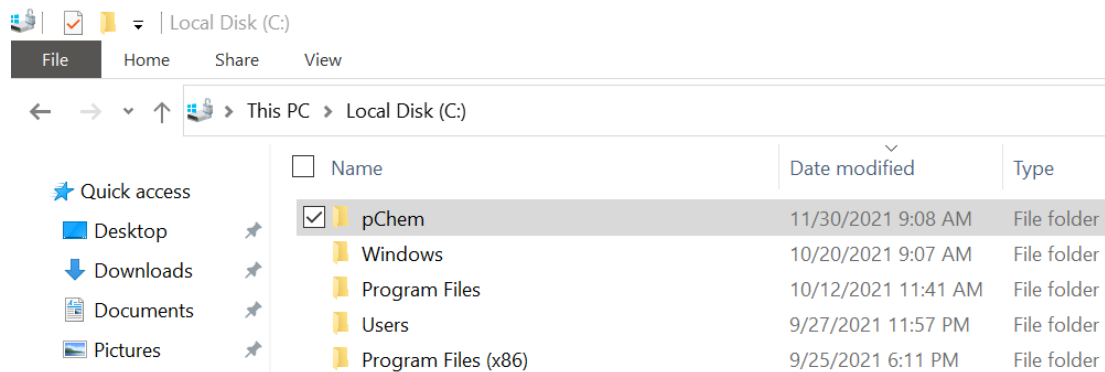
For source code, please refer to [github](#).

For detailed usage, please refer to [user guide](#).


2) click “*click to download*” to get the zipped software.

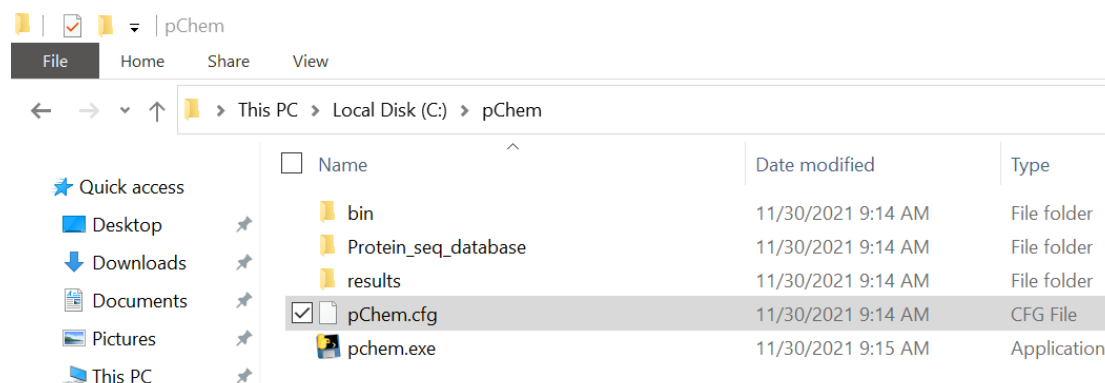


3) Un-zip the “pChem.zip” package into a specified file folder (e.g., Local disk C).



### 3. Configuration

1) Double click *pChem*  *pChem* to open the main folder.



2) Open configuration file "*pChem.cfg*" using a text editor, e.g., Microsoft Notepad or Notepad++ (<https://notepad-plus.en.softonic.com/>).

3) Setting "*pChem.cfg*".

```
1 # If isotope coding is adopted to facilitate the discovery of unknown modifications (True or False)
2 isotope_labeling=True
3
4 # Path to the output file
5 output_path=D:\pchem\pChem_new\results
6
7 # Path to the protein sequence database
8 fasta_path=D:\pChem\pChem_new\Protein_seq_database\Homo_sapiens_uniprot_canonical_20395_entries_20210516.fasta
9
10 # Format of MS data, RAW or MZML
11 msms_type=RAW
12
13 # The number and path of MS data
14 msmsnum=1
15 msmspath1=D:\pchem\pChem_data\QE_Plus_YangJing_FL_ALK_50per_20170531.raw
16
17
18 # Type of MS dissociation method
19 activation_type=HCD-FTMS
20
21 # Usage of open search (True/False), against Unimod, the common modification can be set if not
22 open_flag=False
23 common_modification_number=2...
24 common_modification_list=Carbamidomethyl[C];Oxidation[M];
25
26
27 # Mass tolerance of the mass shift between light isotope and heavy isotope
28 mass_of_diff_diff=6.020132
29
30
31 # Isotopic mass difference within empirically defined tolerance (Da)
32 mass_diff_diff_range=0.005
33
34 # Mass range of unknown modification (Da)
35 min_mass_modification=200
36 max_mass_modification=1000
37
38 # Isotopic pairs of mass shifts with PSMs less than X% of that of overall PSMs were neglected
39 filter_frequency=5
40
41 # If consider the N-side or C-side for amino acid localization (True or False)
42 side_position=True
43
44 # P-value threshold enabling confident amino acid localization
45 p_value_threshold=0.001
46
47 # If report the statistical information (True or False)
48 report_statistics=False
```

### General Note 1:

For the first-time users, custom settings are required for ①-⑤, ⑧ default settings can be adopted for ⑥, ⑦, ⑨-⑭.

### General Note 2:

All parameters (shown in red below) are case sensitive.

### General Note 3:

The blank space should be avoided.

- ① # If isotope coding is adopted to facilitate the discovery of unknown modifications (True or False)

Isotope\_labeling=True

illustration: default

**Note:** Choose 'False', if pChem is adopted to search endogenous modifications from probe-free and/or label-free protein samples,

- ② # Path to the output file

output\_path=C:\pChem\results

**Note:** If the output file folder does not exist, an error will be reported.

- ③ # Path to the protein sequence database

fasta\_path=C:\pChem\Protein\_seq\_database\Homo\_sapiens\_uniprot\_canonical\_20395\_entries\_20210516.fasta

**Note:** The protein \*.fasta database databases of several commonly used species (e.g., *homo sapiens*) are included in the subfolder (named as Protein\_seq\_database) of pChem. Note that the databases of other species can be downloaded from Uniprot as described in **Supporting Protocol 1**.

PC > Local Disk (C:) > pChem > Protein\_seq\_database

<input type="checkbox"/>	Name	Date modified	Type
<input type="checkbox"/>	Arabidopsis_thaliana_uniprot_canonical_16043_entries_20210516.fasta	5/17/2021 12:07 PM	FASTA File
<input type="checkbox"/>	Caenorhabditis_elegans_uniprot_canonical_4226_entries_20210516.fasta	5/17/2021 12:23 PM	FASTA File
<input type="checkbox"/>	Drosophila_melanogaster_uniprot_canonical_3632_entries_20210516.fasta	5/16/2021 11:44 PM	FASTA File
<input type="checkbox"/>	Escherichia_coli_uniprot_canonical_4518_entries_20210516.fasta	5/17/2021 12:15 PM	FASTA File
<input checked="" type="checkbox"/>	Homo_sapiens_uniprot_canonical_20395_entries_20210516.fasta	6/4/2021 9:23 PM	FASTA File
<input type="checkbox"/>	Mus_musculus_uniprot_canonical_17073_entries_20210516.fasta	5/17/2021 12:18 PM	FASTA File
<input type="checkbox"/>	Pseudomonas_syringae_uniprot_canonical_5431_entries_20210516.fasta	7/27/2021 9:56 PM	FASTA File
<input type="checkbox"/>	Rattus_norvegicus_uniprot_canonical_8126_entries_20210516.fasta	5/17/2021 12:22 PM	FASTA File

- ④ # Format of MS data (RAW or MZML)

msmstype=RAW

**Note:** Non-Thermo MS data need to be converted into mzML files before pChem search. The users can refer to **Supporting Protocol 2**.

- ⑤ # The number and path of MS data

msmsnum=N  
msmspath1=X:\XXX\XXX.raw  
msmspath2=X:\XXX\XXX.raw

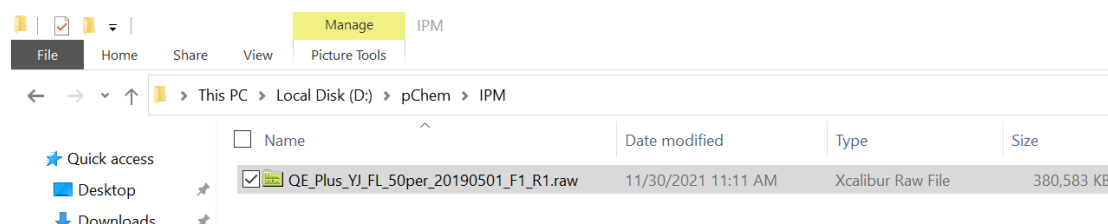
.....  
msmspathN=X:\XXX\XXX.raw

**Note:** The suffix of MS data files MUST be input.

**Example:**

msmsnum=1

msmspath1=D:\pChem\IPM\QE\_Plus\_YJ\_FL\_50per\_20190501\_F1\_R1.raw



⑥ # Type of MS dissociation method

activation\_type=HCD-FTMS

illustration: default

**Note:** 1) Users can adopt this default setting if their MS data is generated by TOF (Time-of-flight) instruments implementing CID (Collision-induced dissociation)-type of fragmentation (e.g., SCIEX 5600, SCIEX 6600, Bruker TimsTOF); 2) pChem v1.0 can NOT support MS data generated under electron-transfer dissociation ETD, electron-transfer/higher-energy collision dissociation EThcD, and the likes.

⑦ # Usage of open search (True/ False) against Unimod, the common modification can be set if not

open\_flag=False

common\_modification\_number=2

common\_modification\_list=Carbamidomethyl[C];Oxidation[M];

illustration: default

**Note:** The names of common modifications should be the same as those appeared in [Unimod](#) database.

⑧ # Mass tolerance of the mass shift between light isotope and heavy isotope

mass\_of\_diff\_diff=6.020132

**Note:** This default value is calculated based on the isotopic mass shift between six heavy and light carbons encoded in probe-derived modifications (PDMs). Users can set any other values based on their different isotope labeling strategies.



**Troubleshooting:** One needs to confirm this value being correctly input.

⑨ # Isotopic mass difference within empirically defined tolerance (Da)

mass\_diff\_diff\_range=0.005

illustration: default

**Troubleshooting:** If the pChem search mis-identified the targeted PDMs or even report nothing, one might want to loose the defined mass tolerance (e.g., 0.01Da).

⑩ # Mass range of unknown modification (Da)

min\_mass\_modification=200

max\_mass\_modification=1000

illustration: default

**Note:** The PDMs generated from the use of bioorthogonal cleavable linkers typically possess masses higher than 200 Da and less than 1000Da.

⑪ # Isotopic pairs of mass shifts with PSMs less than X% of that of overall PDMs were neglected

filter\_frequency=5

illustration: default

**Note:** This parameter can be set as 0 if one wants to retrieve all PDMs including those with just a few PSMs.

⑫ # If consider the N- or C-termini for amino acid localization (True or False)

side\_position=True

illustration: default

⑬ # P-value threshold enabling confident amino acid localization

p\_value\_threshold=0.001


illustration: default

⑭ # if report the statistical information (True or False)

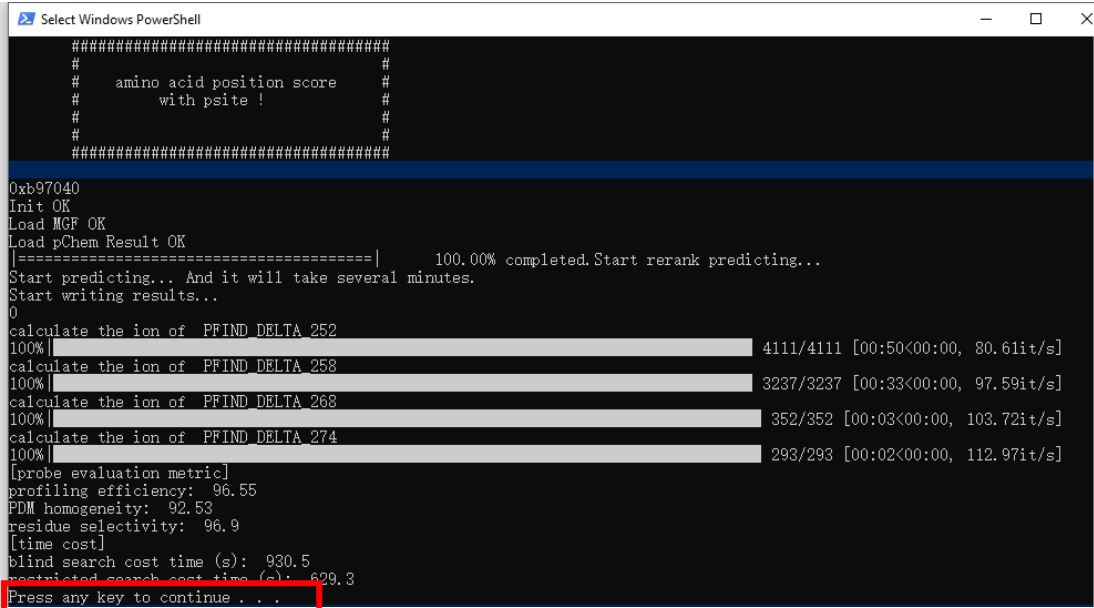
report\_statistics=False

illustration: default

## 4. Run


Once all parameters have been set, double click “*pChem.exe*”  *pchem.exe* to execute the programming. The message “**Please press any key to continue**” means that program runs to completion.

**Note:** pChem search will generate several intermediate files in the main folder. do NOT open those files during program running.

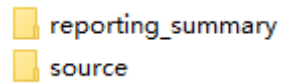


```
#####  
#  
# amino acid position score #  
# with psite ! #  
# #  
# #  
#####  
0xb97040  
Init OK  
Load MCF OK  
Load pChem Result OK  
|=====| 100.00% completed.Start rerank predicting...  
Start predicting... And it will take several minutes.  
Start writing results...  
0  
calculate the ion of PFIND DELTA 252  
100% | 4111/4111 [00:50<00:00, 80.61it/s]  
calculate the ion of PFIND DELTA 253  
100% | 3237/3237 [00:33<00:00, 97.59it/s]  
calculate the ion of PFIND DELTA 263  
100% | 352/352 [00:03<00:00, 103.72it/s]  
calculate the ion of PFIND DELTA 274  
100% | 293/293 [00:02<00:00, 112.97it/s]  
[probe evaluation metric]  
profiling efficiency: 96.55  
PDM homogeneity: 92.53  
residue selectivity: 96.9  
[time cost]  
blind search cost time (s): 930.5  
restricted search cost time (s): 629.3  
Press any key to continue . . .
```

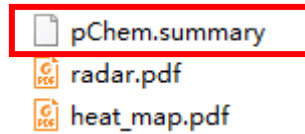
## 5. Output

1) Double click “*results*” file for searching results  *results*.

2) Double click “*reporting summary*”.



3) There are three major output documents.



**Note:** Users are recommended to copy these output documents and paste into another file. Otherwise, they can be covered by those generated from the next search event.

## ① pChem.summary

*pChem.summary* is a tab-delimited text file contains the details of every PDM.

Rank	PDM	Accurate Mass	Top1 Site Probability	Others	#PSM	#PSM L H	DFLs
1	PFIND_DELTA_252	252.122339	C 0.988		13876	7368 6508	
2	PFIND_DELTA_268	268.116411	C 0.487	M(0.291); N-SIDE(0.212);	1578	872 706	302.104737, 301.103528, 320.113184

**PDM:** Probe-derived modifications

**#PSM:** The number of PSMs corresponding to modified peptides identified by targeted search

**#PSM L|H:** The number of PSMs assigning to light and heavy bearing the corresponding PDM, respectively

**Note:** For data from on-isotope-labeled samples, this information will NOT be shown.

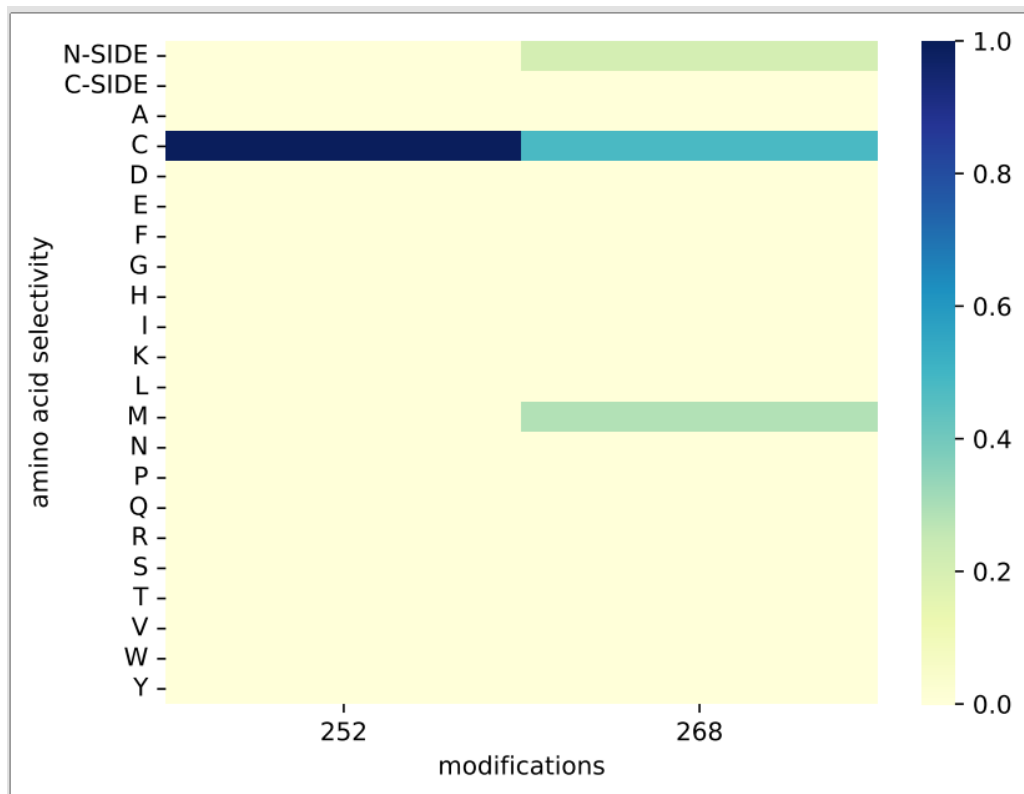
**Top1 site | Top1 Probability:** The amino acid most likely to be modified with the corresponding localization probability.

**Others:** Other amino acid sites that may also be labeled by probes and their corresponding localization probability values

**DFLs:** Diagnostic fragment losses

**Note:** For data from on-isotope-labeled samples, DFLs will NOT be provided.

## ② Heat\_map.pdf



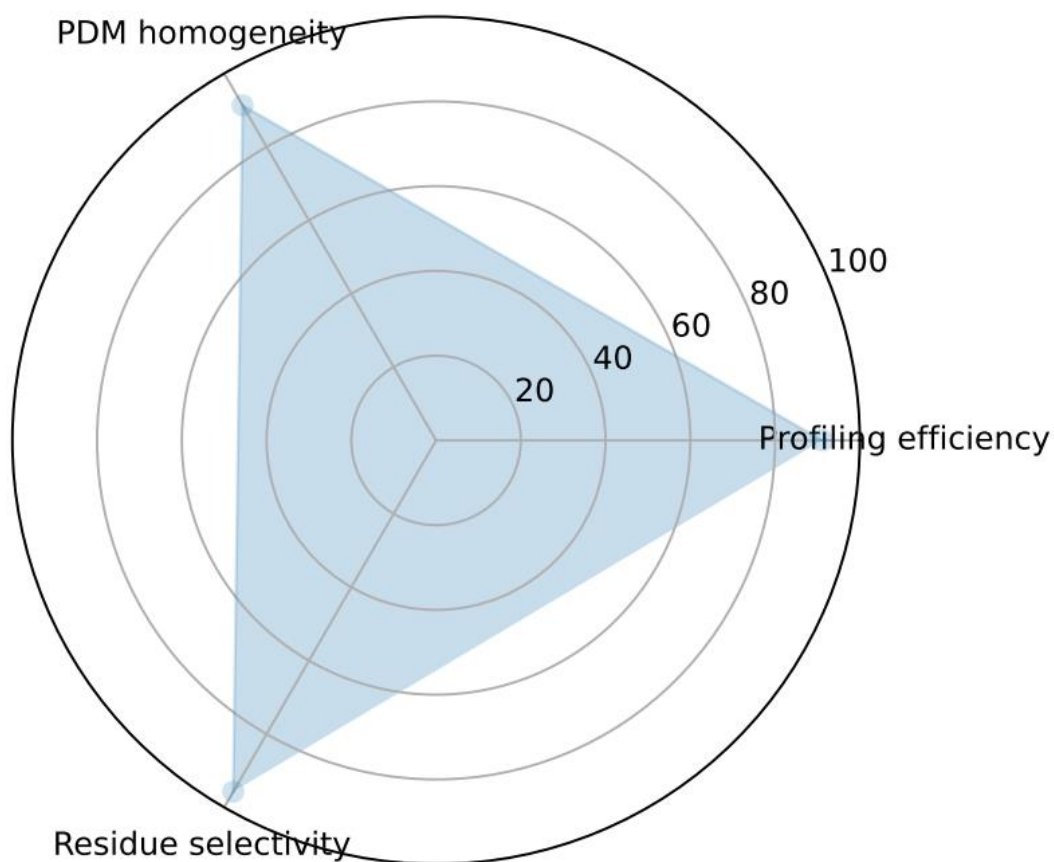
Horizontal coordinate: The  $\Delta$ mass of each PDM

Longitudinal coordinate: The types of amino acids

Color gradient: The localization probability that the modification occurs at each potential site.

**Note:** 1) Those amino acids with p-value higher than 0.001 are considered mis-localized sites. As such, their localization probability values are defined to be null. 2) For data generated from non-isotope-labeled samples, heatmap will NOT be provided.

### ③ Radar.pdf



Radar.pdf contains a radar plot whose radial axes correspond to the three scores as indicated.

**Profiling efficiency (%)** that evaluates whether a probe enables the efficient identification of modified peptides for chemoproteomics.

**Modification homogeneity (%)** that evaluates whether a probe forms a uniform modification.

**Residue selectivity (%)** that evaluates whether a probe selectively targets specific amino acid:

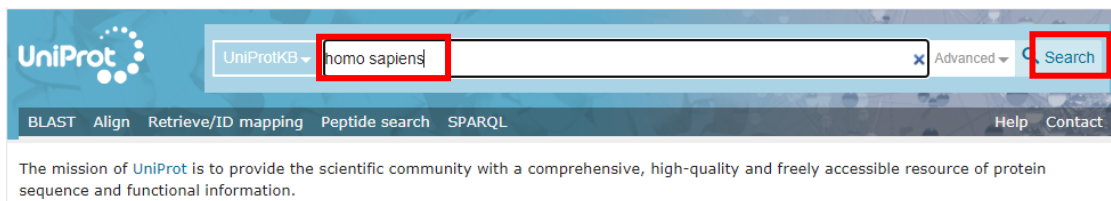
**Note:** 1) *PDM homogeneity* and *Residue selectivity* are calculated based on the blind search results, while *Profiling efficiency* is calculated according to restricted search; 2) For data generated from on-isotope-labeled samples, radar plot will NOT be provided.



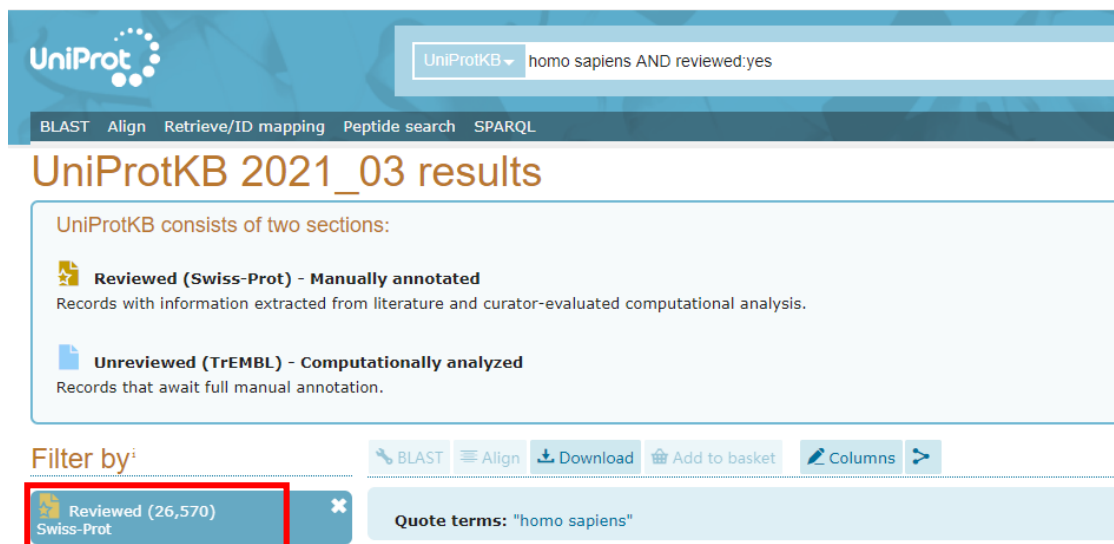
## 6. Supporting protocol 1: Protein sequence database

This protocol is used to download protein \*.fasta files for database search.

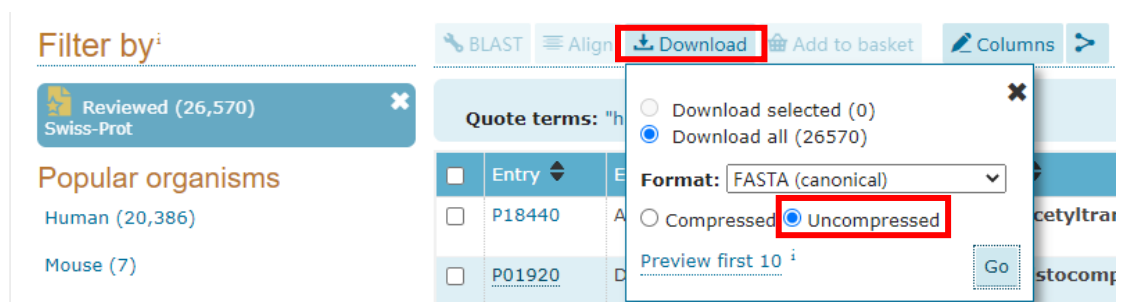
- 1) Open <https://www.uniprot.org/>, enter the Latin name of the species (e.g., *homo sapiens*), then click search.



- 2) Click “Reviewed” (Swiss-Prot).



- 3) Select “Uncompressed”, then Click “Download” and “Go”.



- 4) Get the \*.fasta file.

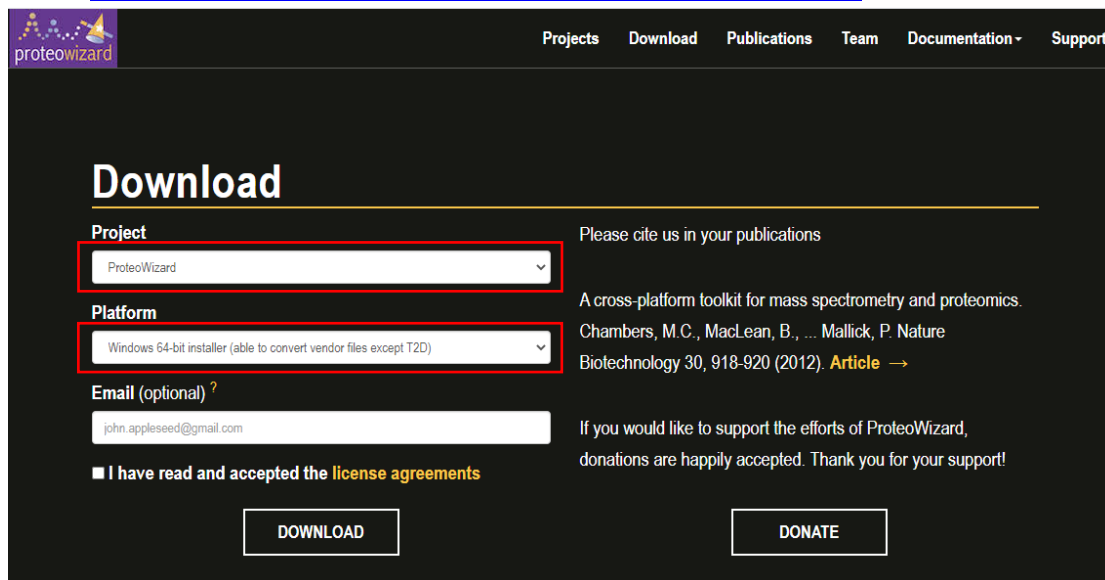




## 7. Supporting protocol 2: MSconvert

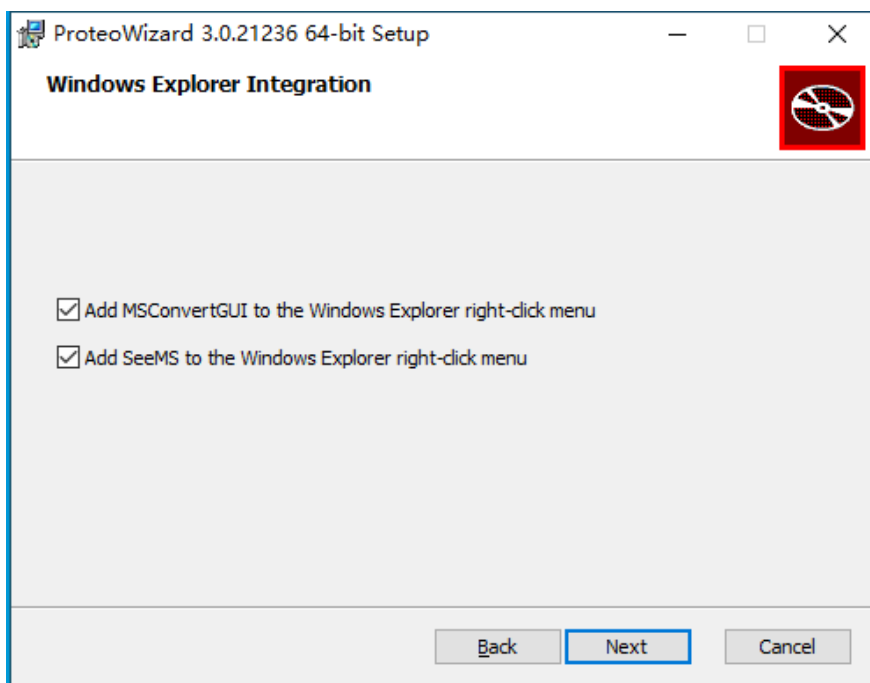
This protocol is used to convert non-Thermo MS data into mzML format files for pChem search.

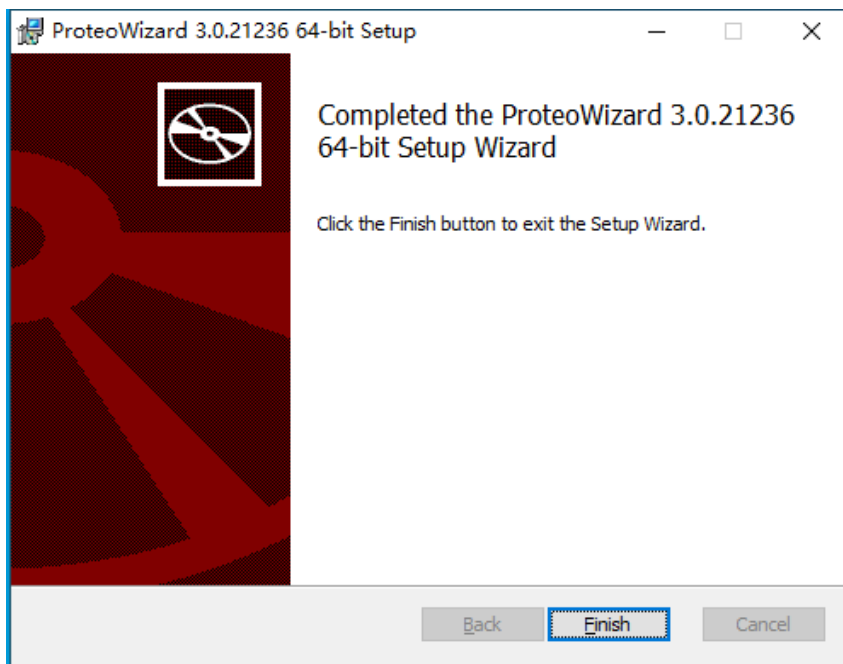
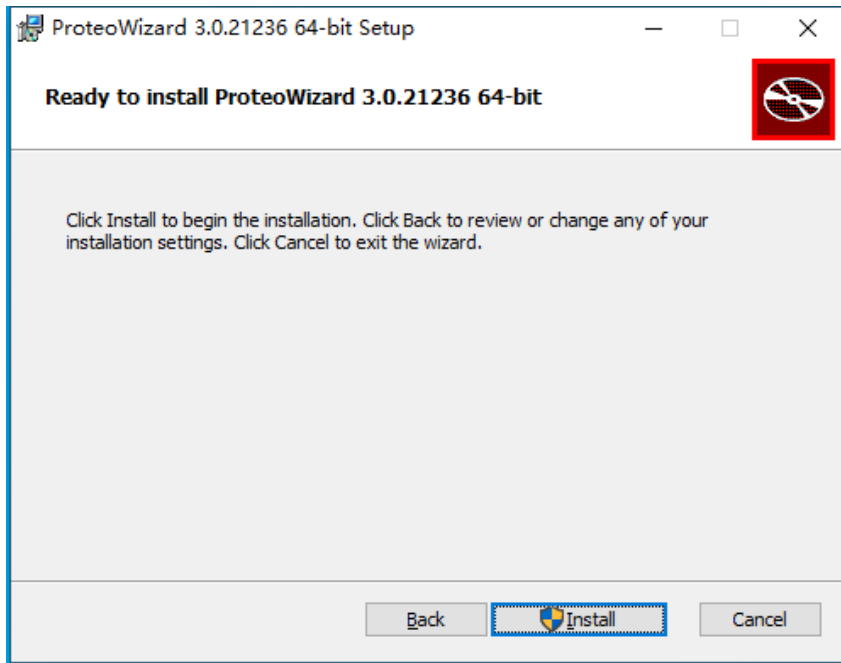
1) Download MSconvertGUI that is embedded in the ProteoWizard platform from: <https://proteowizard.sourceforge.io/download.html> .



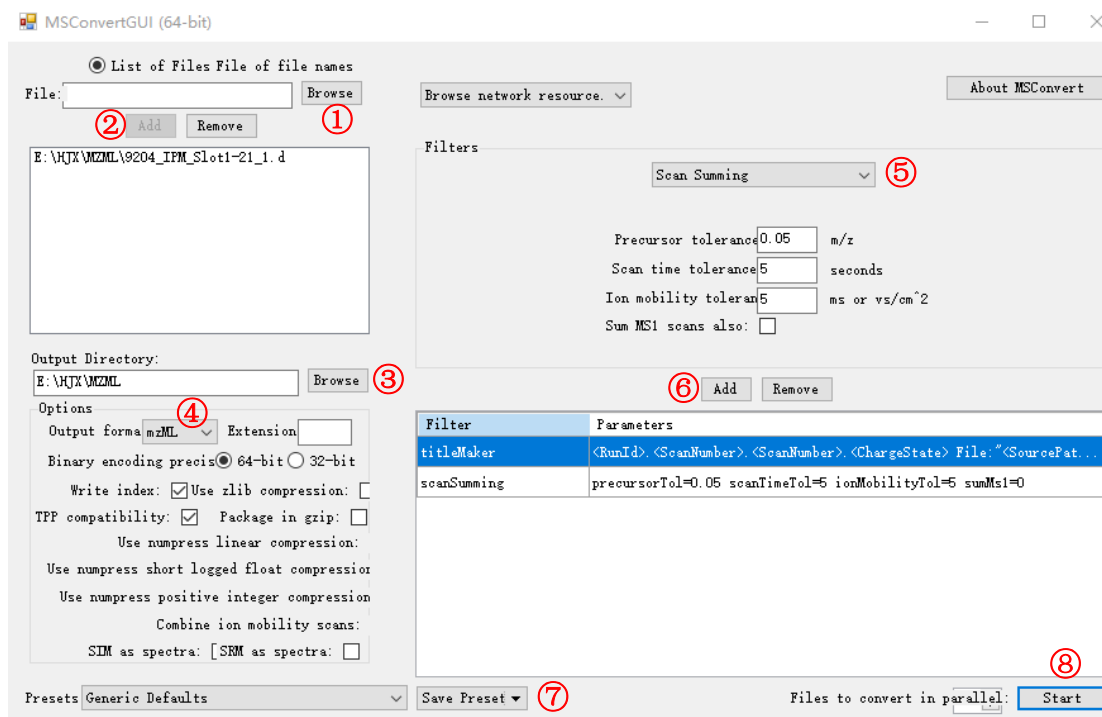
The screenshot shows the ProteoWizard website's download page. The page has a dark background with white text. At the top, there is a navigation bar with links for Projects, Download, Publications, Team, Documentation, and Support. The main heading is "Download". Below the heading, there are two dropdown menus: "Project" (set to "ProteoWizard") and "Platform" (set to "Windows 64-bit installer (able to convert vendor files except T2D)"). Below these is an "Email (optional)" field with the email address "john.appleseed@gmail.com". There is a checkbox labeled "I have read and accepted the license agreements" which is checked. To the right of the form, there is a paragraph of text: "Please cite us in your publications. A cross-platform toolkit for mass spectrometry and proteomics. Chambers, M.C., MacLean, B., ... Mallick, P. Nature Biotechnology 30, 918-920 (2012). Article →". At the bottom of the form, there are two buttons: "DOWNLOAD" and "DONATE".

2) Install ProteoWizard according to the following instruction.





### 3) Open MSConvertGUI



①-② Browse and add MS data (e.g., \*.d, \*.WIFF files)

③ Define output route

④ Choose \*.mzML as the output data format

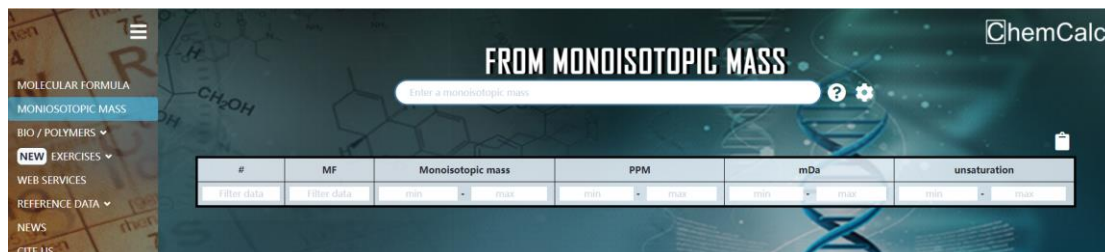
⑤ -⑥ Define parameters for Scan Summing

⑥ -⑧ Save and run

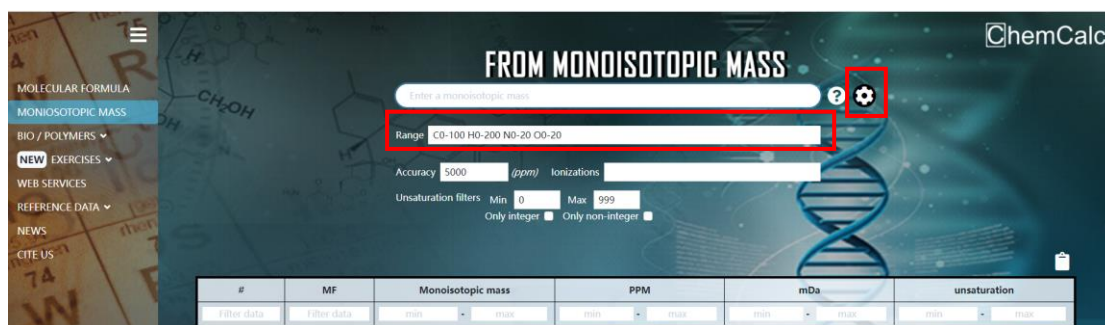
## 8. Supporting protocol 3: ChemCalc

This protocol is used to estimate candidate molecular formulas from the pChem-determined accurate masses.

1) Open <https://www.chemcalc.org/mf-finder>.



2) Click , check the element composition.



3) Input the monoisotopic mass of each PDM shown in *pChem.summary* file. The candidate molecular formulas will immediately appear below.

Rank	PDM	Accurate Mass	Top1 Site Probability	Others	#PSM	#PSM L H	DFLs
1	PFIND_DELTA_252	252.122339	C 0.988		13876	7368 6508	

