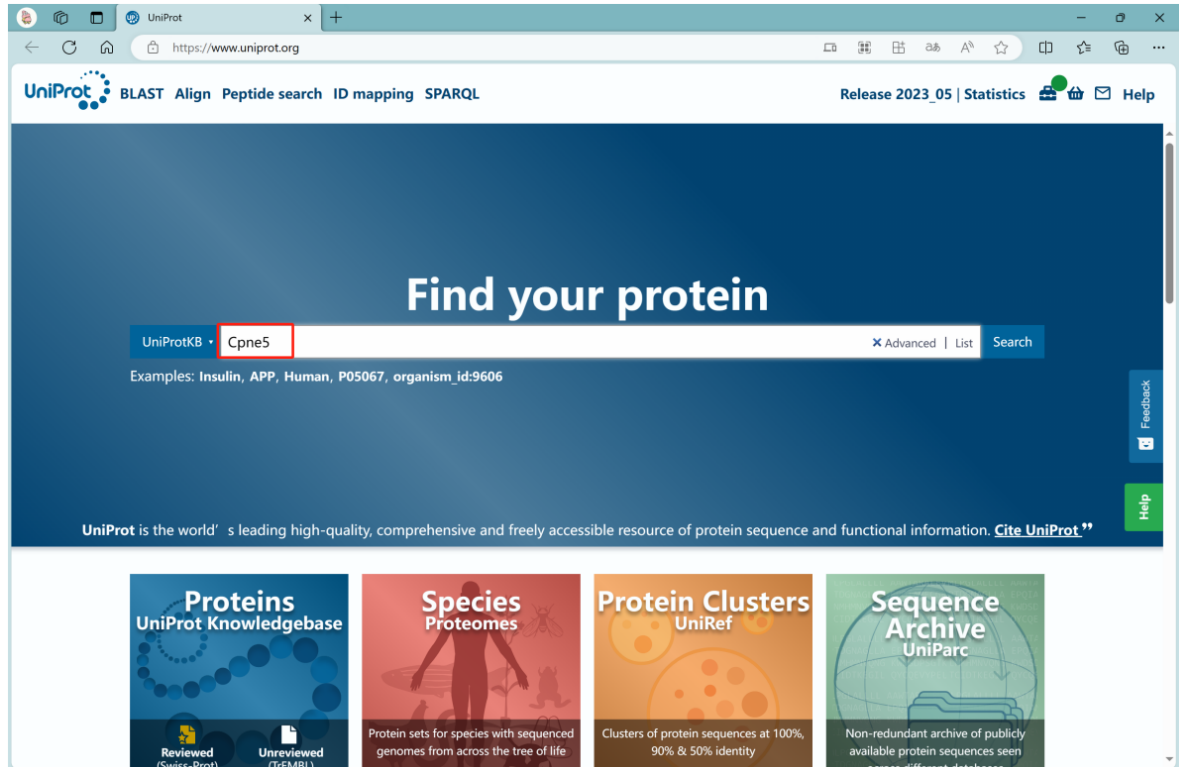


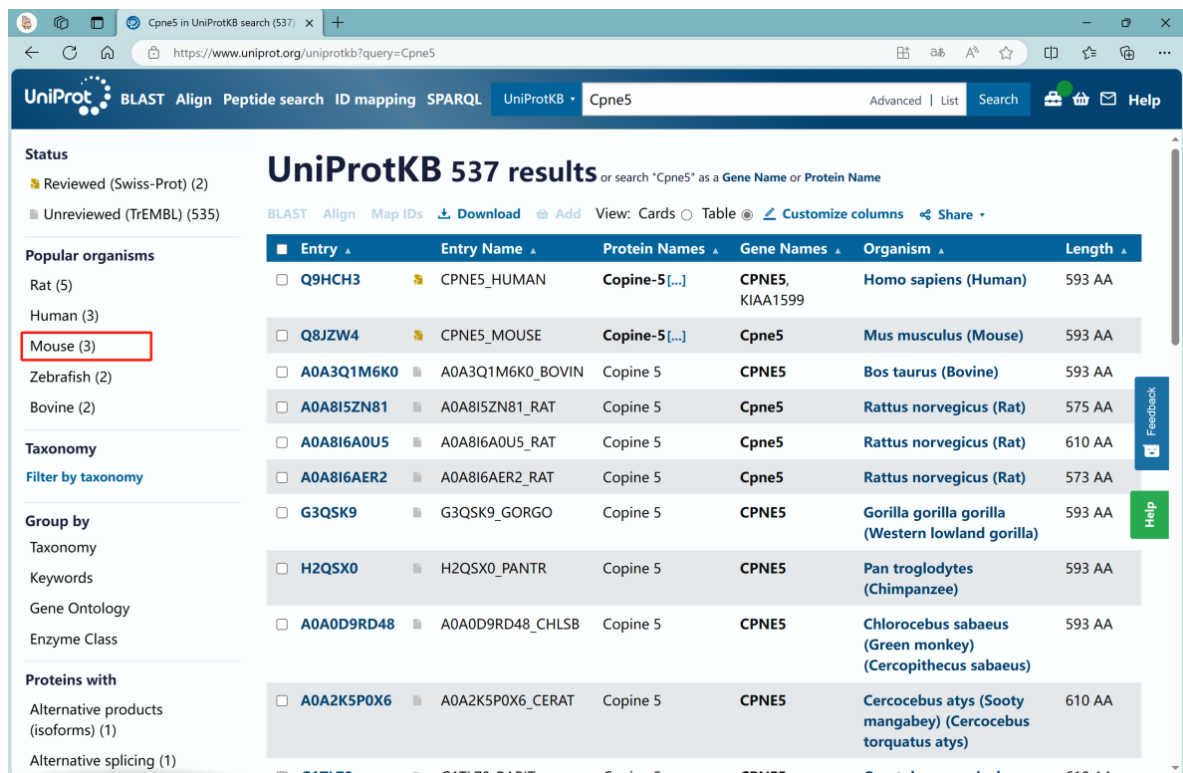
This part of the tutorial will guide you how to obtain the ESM-2 embedding of the gene.

Query

Go to [UniProt](https://www.uniprot.org), input the gene name in the search box (**Cpne5(Mouse)** is an example here), and press **Search**:



Then **select Mouse(or else you need)** in the left filter box:



Make sure the red box is Cpne5, **click the yellow box** (select the first one if the Gene Names are the same), then **click Download** in the green box:

UniProtKB 3 results

BLAST Align Map ID **Download** Add View: Cards Table Customize columns Share 1 row selected out of 3

Entry	Entry Name	Protein Names	Gene Names	Organism	Length
<input checked="" type="checkbox"/> Q8JZW4	CPNE5_MOUSE	Copine-5[...]	Cpne5	Mus musculus (Mouse)	593 AA
<input type="checkbox"/> E9Q1V2	E9Q1V2_MOUSE	Copine V	Cpne5	Mus musculus (Mouse)	307 AA
<input type="checkbox"/> Q8BVA1	Q8BVA1_MOUSE	C2 domain-containing protein	Cpne5	Mus musculus (Mouse)	153 AA

After determining that the red box is the FASTA format for downloading the selected 1 gene (this is the default), **click the yellow box** to download:

Download

☒ Download selected (1)
☐ Download all (3)

Format
 FASTA (canonical)

Compressedⁱ
☒ Yes
☐ No

Generate URL for API Preview 1 Cancel **Download**

When you download and unzip it, you get a file with the gene name in red box and the protein expressions in yellow box:

```
uniprotkb_accession_Q8JZW4_2024_01_06.fasta - 记事本
文件(E) 编辑(E) 格式(O) 查看(V) 帮助(H)
>sp|Q8JZW4|CPNE5_MOUSE Copine-5 OS=Mus musculus OX=10090 GN=Cpne5 PE=1 SV=1
MEQPEDMASLSEFDSLGSIPATKVEITVSCRNLLDKDMFSKSDPLCVMYTQGMENKQWR
EFGRTVIDNTLNPDFVRKFIVDYFFEEKQNLRFDLVDVDSKSPDLKSHDFLGQAFCTLG
EIVGSSGSRLEKPLTIGTFLNSRTGKPMFPAVSNGGVPGKKCGTIILSAEELSNCRDVAT
MQFCANKLDKKDFFGKSDPFLVFYRSNEDGFTTICKTEVMKNTLNPVWQTFSPVRALC
NGDYDRTIKVEVYDWDNRDGSDFIGFTTSYRELARGQSQFNIYEVINPKKKMKKKKYVN
SGTVTLLSFAVESESTFLDYIKGGTQINFTVAIDFTASNGNPSQSTSLHYMSPYQLNAYA
LALTAVGEIIGHYDSDKMFALGFAGLPPDGRVSHEFPLNGNQENPSCCGIDGILEAYH
SSLRTVQLYGPTNFAPVVTHVARNAAAVQDGSQYSVLLIITDGVISDMAQTKEAIVNAAK
LPMSIIIVGVGQAEFDAMVELDGDVRISSRGKLAERDQVQFVFRDYVDRTGNHVLMSA
RLARDVLAIEIPDQLVSYMKAAQIRPRPPPAAPQSPQSPAHSPGSPVHTHI
```

Extract the gene name and protein expression, the following files are obtained:

```
uniprotkb_accession_Q8JZW4_2024_04_16.fasta - 记事本
文件(E) 编辑(E) 格式(O) 查看(V) 帮助(H)
>Cpne5
MEQPEDMASLSEFDSLGSIPATKVEITVSCRNLLDKDMFSKSDPLCVMYTQGMENKQWREFGRTEVIDNTLNPDFVRKFIVE
```

ESM-2 Embedding

Installed the ESM-2 model (you can refer to [ESM-2](#)), then **enter the following command**:

```
cd esm-main/

python scripts/extract.py esm2_t36_3B_UR50D
uniprotkb_accession_Q8JZW4_2024_04_16.fasta examples/data/some_proteins_emb_esm2
--repr_layers 36 --include mean_per_tok
```

Where `uniprotkb_accession_Q8JZW4_2024_04_16.fasta` is the gene name and protein expression file processed above.

After running successfully, the gene embedding file (in this case, `Cpne5.pt`) is generated in the `esm-main/examples/data/some_proteins_emb_esm2/` directory.

After converting all the genes of your dataset to some `*.pt` files, **put them in a folder** (such as `./pt/`) and **run the following python code**:

```
import os
import torch
import pickle
import pandas as pd

df = pd.DataFrame()

for path, dir_lst, file_lst in os.walk(r'./pt'):
    for file_name in file_lst:
        data = torch.load(open(os.path.join(path, file_name), 'rb'))
        print(data['label'])
        print(data['representations'][36][-1])
        df.insert(df.shape[1], data['label'], data['representations'][36]
[-1].numpy())
```

```
pickle.dump(df, open('emb.pkl', 'wb'))
```

Then the ESM-2 embedding file `emb.pkl` of the genes can be generated.