

Aufgabe 1:

```
1  # version 3.8.0
2  import os
3  import urllib.request
4  import random
5  import gzip
6  import re
7  import pprint
8  import sqlite3
9
10 pp = pprint.PrettyPrinter(indent=2)
11
12 filePath = "./ftpSourceView.txt";
13
14 fileContent = open(filePath).readlines();
15
16 sourceUrl = fileContent[0].split()[1]
17 print("URL: " + sourceUrl)
18
19 unzippedFiles = []
20
21
22 def get_data(nr_files):
23     pp.pprint(">>>> STARTING DOWNLOAD <<<<<")
24     toReturn = []
25     dead_files = 0
26     for i in range(nr_files):
27         data = {}
28         fullUrl = sourceUrl + fileContent[2:][random.randint(0, len(fileContent) - 2)].split('')[1]
29         try:
30             content = gzip.GzipFile(fileobj=urllib.request.urlopen(fullUrl)).readlines()
31             matched_header = re.search(r"HEADER\s*([\w\s-]*)\d{2}-\w{3}-\d{2}\s*(\d*\w*)", content[0])
32             data["header"] = matched_header.group(1).strip()
33             data["pdbname"] = matched_header.group(2)
34             data["X"] = 0
35             data["Y"] = 0
36             data["Z"] = 0
37             data["title"] = []
38             for line_src in content[1:]:
39                 line = line_src.decode()
40                 match = re.match(r"TITLE\s\d?\s?(.*)", line)
41                 if match != None:
42                     data["title"].append(match.group(1).strip())
43                 elif line.startswith("ATOM"):
```

```
44         split = line.split();
45         data["X"] += float(split[6])
46         data["Y"] += float(split[7])
47         data["Z"] += float(split[8])
48     else:
49         match = re.match(r".*CHAIN:\s(.*);", line)
50         if match != None:
51             data["chain"] = match.group(1)
52         data["title"] = " ".join(data["title"])
53         pp.pprint(str(((i + 1) / nr_files) * 100) + "% completed")
54         toReturn.append(data)
55     except:
56         pp.pprint("got a faulty file, continuing ...")
57         dead_files += 1
58 pp.pprint(">>>>> DONE <<<<<")
59 pp.pprint("!! Encountered " + str(dead_files) + " files.")
60 return toReturn
61
62 def insert_data(data, database_file):
63     pp.pprint(">>>>> STARTING DB INSERT <<<<<")
64     con = sqlite3.connect(database_file)
65     cur = con.cursor()
66     for i, element in enumerate(data):
67         cur.execute('insert into proteins (' +
68                     'header, title, pbid, chain, X, Y, Z' +
69                     ') VALUES (?, ?, ?, ?, ?, ?, ?)',
70                     [element["header"], element["title"], element["pbid"],
71                     element["chain"], element["X"], element["Y"], element["Z"]])
72     pp.pprint(str(((i + 1) / len(data)) * 100) + "% completed")
73     con.commit()
74     con.close()
75     pp.pprint(">>>>> DONE WITH DB INSERT <<<<<")
76
77 def setup_db(database_file):
78     con = sqlite3.connect(database_file)
79
80     cur = con.cursor()
81
82     cur.execute('create table if not exists proteins (' +
83                 'id integer primary key autoincrement,' +
84                 'header text,' +
85                 'title text,' +
86                 'pbid text,' +
87                 'chain text,' +
88                 'X float,' +
89                 'Y float,' +
```

```
90         'Z float' +  
91         ')')  
92     con.commit()  
93     con.close()  
94  
95     database_file = "./protein.db"  
96  
97     setup_db(database_file)  
98     insert_data(get_data(10), database_file)
```

SQL script:

```
1  # version 3.8.0
2  import sqlite3
3  import pprint
4
5  pp = pprint.PrettyPrinter(indent=2)
6
7  con = sqlite3.connect('./protein.db')
8
9  cur = con.cursor()
10
11 result = cur.execute("select * from proteins where pbid like ('%5%')")
12
13 for line in result.fetchall():
14     pp.pprint(line)
15
16 con.commit()
17 con.close()
```

Aufgabe 2:

```
1  # version 3.8.0
2  import urllib.request
3  import re
4  import pprint
5
6  PP = pprint.PrettyPrinter(indent=1)
7
8  BASE_URL = "https://www.icd-code.de/icd/code/"
9  STARTING_URL = "ICD-10-GM.html"
10 LIST_OF_IDS = ["H", "J"]
11
12
13 def build_tree(url, ids, base):
14     tree = {}
15     file = urllib.request.urlopen(url)
16     for j in file.readlines():
17         for i in ids:
18             identification = re.search(r"<a href=\"(" + i + r"\d\d-" + i +
19                                     r"\d\d.html)\">>(\w\d\d-\w\d\d)</a></td><td>(.*?)<a",
20                                     str(j))
21             if identification is None:
22                 continue
23             key = identification.group(2)
24             tree[key] = {
25                 "link": base + identification.group(1),
26                 "detail": base + identification.group(3),
27                 "more": parse_second_level(i, base, base + identification.group(1))
28             }
29     PP.pprint(tree)
30
31 # def parse_first_level(url, ids, base, tree):
32
33
34 def parse_second_level(cur_id, base, link):
35     tree = {}
36     for line in urllib.request.urlopen(link).readlines():
37         matches = re.findall(r"<a href=\"(" + cur_id +
38                             r"\d\d-\w\d\d.html)\">>(\w\d\d-\w\d\d)</a>(.*?)(<br/>|</td>)",
39                             str(line))
40         if matches is None:
41             continue
42         for match in matches:
43             key = match[1]
```

```

44     sublink = base + match[0]
45     tree[key] = {
46         "link": sublink,
47         "detail": match[2].strip(),
48         "more": parse_third_level(cur_id, base, sublink)
49     }
50     return tree
51
52
53 def parse_third_level(cur_id, base, link):
54     tree = {}
55     for line in urllib.request.urlopen(link).readlines():
56         matches = re.findall(r"<a href=\"(" + cur_id +
57                               r"\d\d(?:\.\d)??.html)\"><(" + cur_id +
58                               r"\d\d(?:\.\d)??)</a>(.*?)(<br/>|</td>)", str(line))
59         if matches is None:
60             continue
61         for match in matches:
62             key = match[1]
63             sublink = base + match[0]
64             tree[key] = {
65                 "link": sublink,
66                 "detail": match[2].strip(),
67                 "more": parse_fourth_level(cur_id, base, sublink)
68             }
69     return tree
70
71
72 def parse_fourth_level(cur_id, base, link):
73     tree = {}
74     for line in urllib.request.urlopen(link).readlines():
75         matches = re.findall(r"<tr><td valign=\"top\"><div class=\"code_bottom\">(" + cur_id +
76                               r"\d\d\d\.\d\d?(?:-?\\*)?)</div></td><td colspan=\"2\">(.*?)</td></tr>",
77                               str(line))
78         if matches is None:
79             continue
80         for match in matches:
81             key = match[0]
82             tree[key] = match[1].strip()
83     return tree
84
85
86 build_tree(BASE_URL + STARTING_URL, LIST_OF_IDS, BASE_URL)

```

Aufgabe 2:

```
1  # version 3.8.0
2  import re
3  import pprint
4
5  pp = pprint.PrettyPrinter(indent=2)
6
7  emailFile = open("./imgtSpeciesGeneInformation.txt", "r")
8
9  # pattern = re.compile(r"((\w*\s?)*);((\w*\d*-*?)*);")
10
11  res_dict = {}
12  avail_gen = []
13  line_nr = 0
14
15  for i in emailFile.readlines()[1:]:
16      mo = re.search(r"((\w*\s*)*);((\w|\d|-|/)*);.*", i)
17      line_nr += 1
18      if mo != None:
19          if len(mo.group(3)) < 4:
20              continue
21          if mo.group(1) not in res_dict:
22              res_dict[mo.group(1)] = {}
23          if mo.group(3)[3] not in res_dict[mo.group(1)]:
24              res_dict[mo.group(1)][mo.group(3)[3]] = 0
25          res_dict[mo.group(1)][mo.group(3)[3]] += 1
26          avail_gen.append(mo.group(3)[3])
27
28  avail_gen = set(avail_gen)
29
30  for g in avail_gen:
31      for k in res_dict:
32          if g not in res_dict[k]:
33              res_dict[k][g] = 0
34
35  # TODO: final calculation missing
36
37  pp.pprint(res_dict)
38
39  emailFile.close()
```
