

**Софийски университет „Св. Климент Охридски”,
Факултет по математика и информатика**



Курсов проект

**Предмет: Линукс и езици за програмиране в
биоинформатиката**

Зимен семестър, 2020/2021

**Преподавател: доц. д-р Д. Василев
Преподавател: Илиян Михайлов**

**Изготвил: Моника Борисова,
ф.№ 26143, БМИ**

Февруари, 2021

София

Цел на проекта:

Да се разработи RESTFul приложение на Python (по желание може и друг език), което да оперира върху данни от ensembl.

Задачи:

1. Да приема id на gene и да връща като резултат цялата секвенция на гена и всички екзони като отделни свойства.

Извеждане на секвенцията на гена в json формат:

```
27 # getting sequence and exons, Get request by default
28 @app.route('/v1/sequence/gene/id/<id>/')
29 def get_sequence(id):
30     # in json format
31     sequence = requests.get(server + gene_id + id,
32                             headers={"Content-Type": "application/json"}).json()
33     exons = requests.get(server + exon_ext + id + '?expand=1',
34                           headers={"Content-Type": "application/json"}).json()
35
36     # from the dictionary sequence, getting the value behind key seq
37     seq = sequence["seq"]
38
```

Seq е String, който съдържа стойността на ключа seq в речника sequence.

За да вземем транскриптите и секвенциите в резултата, append-ваме речници, и резултата ни става списък, който се състои от малки речници с дадените стойности и ключове.

```
6 app = Flask(__name__)
7
8 server = "https://rest.ensembl.org"
9 gene_id = "/sequence/id/"
10 exon_ext = "/lookup/id/"
11
12 # Transcripts and Exon function
13 def get_gene(transcripts, sequence):
14     result = list()
15
16     for transcript in transcripts:
17         for exon in transcript['Exon']:
18             result.append(
19                 {"start": exon["start"],
20                  "end": exon["end"],
21                  "id": exon["id"],
22                  "seq": sequence[exon["start"]:exon["end"]]
23             })
24
25     return result
```

Output:

[illegible]

2. Да приема id на gene и да връща като резултат секвенцията с пресметнато процентното GC съдържание и разменени аминокиселини аденин с тимин в конкретния пример.

```
gcContent = request.args.get("gc_content")
swap = request.args.get("swap")
contentType = request.args.get("content-type")

# if we have ?gc_content=true&swap=A:T
if (gcContent == "true" or swap is not None):
    result = {"seq": seq}

    # check if gcContent is true and return calculated G and C content. The percentage is calculated against the full length. (0-100)
    if (gcContent == "true"):
        result["gc_content"] = GC(seq)

    # check if swap is given and split by : for (A:T), swap the amino acids A with T
    if (swap is not None):
        swap = swap.split(":")
        result["swap_sequence"] = swapFunction(seq, swap[0], swap[1])

    return result

# fasta, x-fast, multi-fast
elif contentType:
    return sequenceByContentType(id, sequence, contentType)

# default whole sequence and all the exons
else:
    return gene = get_gene(exons.get('Transcript'), str(sequence['seq']))
    return {'seq': sequence, 'exons': returnExons}
```

От URL искаме да вземем това, което стои след `?(args)`, в случая:

"gc_content", "swap", "content-type", через getter.

След което проверяваме дали gcContent и swar не са None, конкретно искаме gcContent да ни е “true”, за да върнем секвенцията, с два if-а пресмятаме процентното GC съдържание, пресмята се с готов метод от модула Biopython и разменяме базите Аденин с Тимин със следната функция:

```
# swap function for A and T
def swapFunction(seq, n1, n2):
    seq = seq.replace(n1, 'S')
    seq = seq.replace(n2, n1)
    seq = seq.replace('S', n2)

    return seq
```

[illegible]

Функция за различните формати с 2 функции за форматите - x-fasta, multi-fasta:

```
# different formats in query
def sequenceByContentType(id, sequence, contentType):
    if contentType == "fasta":
        return {"seq": sequence, "id": ">" + id}
    elif contentType == "x-fasta":
        return xFasta(id, sequence)
    elif contentType == "multi-fasta":
        return multiFasta(id, sequence)

def xFasta(id, sequence):
    description = ">" + sequence["id"]
    return description + "\n" + sequence["seq"]

def multiFasta(id, sequence):
    result = ""
    for seq in sequence:
        desc = ">" + sequence["id"] + "." + \
            str(sequence["version"]) + " " + sequence["desc"]
        result += desc + "\n" + seq

    return result
```

```
* Serving Flask app "project" (lazy loading)
* Environment: production
WARNING: This is a development server. Do not use it in a production deployment.
Use a production WSGI server instead.
* Debug mode: off
* Running on http://127.0.0.1:5000/ (Press CTRL+C to quit)
127.0.0.1 - - [05/Feb/2021 12:48:55] "GET /v1/sequence/gene/id/ENSG00000157764/ HTTP/1.1" 200 -
127.0.0.1 - - [05/Feb/2021 12:48:57] "GET /favicon.ico HTTP/1.1" 404 -
127.0.0.1 - - [05/Feb/2021 12:50:24] "GET /v1/sequence/gene/id/ENSG00000157764/?gc_content=true&swap=A:T HTTP/1.1" 200 -
127.0.0.1 - - [05/Feb/2021 12:51:18] "GET /v1/sequence/gene/id/ENSG00000157764/?content-type=fasta HTTP/1.1" 200 -
127.0.0.1 - - [05/Feb/2021 12:52:17] "GET /v1/sequence/gene/id/ENSG00000157764/?content-type=x-fasta HTTP/1.1" 200 -
127.0.0.1 - - [05/Feb/2021 12:52:53] "GET /v1/sequence/gene/id/ENSG00000157764/?content-type=multi-fasta HTTP/1.1" 200 -
```