

protein seqfasta dictionary function

October 11, 2020

```
[35]: x = ">gi|27544883|ref|NP_775396.1| putative ribosomal protein S3_1
      ↪(mitochondrion) [Lecanicillium muscarium]
      MEKNIKLNILLSKVPQKITLNNKIADIRKTKYLPFSKEWKDTIYSYNKNIMKNIPSHHLNINKIIQSYFNLFFSTNKN
      NQNKRFISFGKYITMKRRRNLLRKIYVSNPVFKFSYFSAFITLFSLSREKVFYKKNILKRIKKYKVIKCNVIYFYIKI
      WIYKIFLAKQYYKNKNISYFLLKNKNKFIQYKLKYLKFLLLKNLYLKRVSWSKIKNFIKRHLIFLRKYELLYSLNQLKF
      NKLTLLNKLSTLLNKLKILGKKIEYNIINLKSIIFNSDLFTQAITLKFQKRKSFNYKKNILSILGRVNFYFKDVVALAYNE
      STNYILNKYKDGKILSYINNHQNLNDFINKIHNNATNKNIHKEIFNSIQYKNIEGIRIETNGRLTKRYRADRAVHYRKWKG
      GLQKTSLSNSTLFRGNVNPNISSIVNTRRVGSFAIRGWISGK
      >gi|27544884|ref|NP_775397.1| NADH dehydrogenase subunit 2 (mitochondrion)_1
      ↪[Lecanicillium muscarium]
      MIIISILSLLLSNAVNLRDVSILYNRIAILILLYCIVHDYTSLTVVTKGIGLHGGLLLINNLTQIFHIFVFIVTIFILT
      LTSFYPRKVWVSEYSSIKDLEPVDPLFNKFIYNTKIINKMGEHLKIIIEYPLILLFIVRPGAIFLMSSNDLITIFLSIEL
      QSYGLYILSTVYRNSELSTTGGLIYFLLGGLSSCFILLGTALLYAKSGRTSLESYIITSISDIHSSTNDLWYSPNYISL
      SLVIFTVGFLFKISAAPFHFVSPDVDAIPTIVTTFVALIAKISILILLQLVYYTNADITMNNWTFILLISSLFTLVIG
      TVVGLTQFRIKRLFAYSTISQIGFMLLALSISSIESTQAPIFYLTQYIIINLNAFIILLAITYSLYFYTNYNKEHKDLLD
      KTNSPIQLITQLKGYFFINPVLALNLTITIFSFSQIPLLVCFCKQIVLSASLNQNLYFITLIAILTNFIERIYYLTIK
      KIFFNKSDYKINTLLCNFKLKINIYNNMNSINSVEYNYKNIFLSIPISFIISILTLTILFFLFINKKWLIMITIFVQLIF
      NS
      >gi|27544885|ref|NP_775398.1| NADH dehydrogenase subunit 3 (mitochondrion)_1
      ↪[Lecanicillium muscarium]
      MSGTTFLFIFVCVIAILFLALNFILAPHNPYQEKYSIFECGFHSFLGQNRTQFGFKFFIFSLVYLLLDLDLEILVIYPYG
      LSSYENGIIYGLIIVLLFIGIITAGFVFELGKGALKIDSRQSYNYFNVQSTKNFINTFFENK"""
```

```
[36]: def protein_dict(seq):
      pro_dict = {}
      seq_split = seq.split('>')
      seq_split.pop(0)
      for i in range(len(seq_split)):
          for pro in seq_split:
              pro_dict[seq_split[i].split('\n')[0]] = ''.join(seq_split[i].
              ↪split('\n')[1 : len(pro) + 1])
      return(pro_dict)
```

```
[37]: protein_dict(x)
```

```
[37]: {'gi|27544883|ref|NP_775396.1| putative ribosomal protein S3 (mitochondrion)
      [Lecanicillium muscarium]': 'MEKNIKLNILLSKVPQKITLNNKIADIRKTKYLPFSKEWKDTIYSYNKN
```

```

IMKNIPSHHLNINKIIQSYFNLFFSTNKNQNKRFISFGKYITMKRRRNLLRKIYVSNPVFKFSYFSAFITLFSLSREKV
FYKKNILKRIKKYKVIKCNVIYFYIKIWKIYKIFLAKQYYKNKNISYFLLKNKNKFIQYKLYLSKFLLLKNLYLKRWW
SKIIKNFIKRHLIFLRKYELLYSLNQLKFNKLTLLNKLSSLLNKILGKKIEYNIINLKSIIFNSDLFTQAITLKFQKRKS
FNYKKNILSILGRVNFYFKDVVALAYNESTNYILNKYKDGKILSYINNHQNLDNFINKIHNAATKNIHKEIFNSIQYKN
IEGIRIETNGRLTKRYRADRAVHYRKWKGGGLQKTSLSNSTLFRGNVNPNISYSIVNNTRRVGSFAIRGWISGK',
'gi|27544884|ref|NP_775397.1| NADH dehydrogenase subunit 2 (mitochondrion)
[Lecanicillium muscarium]': 'MIIISILSLLLSNAVNLRDVSILYNRIAILILLYCIVHDYTSLTVVTGKI
GLHGGLLLINNLTIQIFHIFVIVTIFILTTSFYPRKVVWSEYSSIKDLEPVDPLFNKFIYNTKIINKMGEHLKIIIEYP
LILLFIVRPGAIFLMSSNDLITIFLSIELQSYGLYILSTVYRNSELSTTGGLIYFLLGGLSSCFILLGTALLYAKSGRTS
LESLYIITSISDIHSSNDLWYSPNYISLSLVIFTVGFLFKISAAPFHFWSPDVYDAIPTIVTTFVALIAKISILILLQ
LVYYTNADITMNNWTFILLISSLFTLVIGTVVGLTQFRIKRLFAYSTISQIGFMLLALSISSIESTQAPIFYLTQYIIIN
LNAFIILLAITYSLYFYTNYNKEHKDLLDKTNSPIQLITQLKGYYFINPVLALNLTITIFSFSQIPPLLVCFCQKIVLSA
SLNQNLFITLIAILTNFIERIYYLTIKKIFFNKSDYKINTLLCNFKLKINIYNNNSINSVEYNYKNIFLSIPISFII
SILTTLTILFFLFINKKWLIMITIFVQLIFNS',
'gi|27544885|ref|NP_775398.1| NADH dehydrogenase subunit 3 (mitochondrion)
[Lecanicillium muscarium]': 'MSGTTFLFIFVCVIAILFLALNFILAPHNPYQEKYSIFECGFHSFLGQNRT
QFGFKFFIFSLVYLLLDLEILVIYPYGLSSYENGIYGLIIVLLFIGIITAGFVFELGKGALKIDSRQSYNYFNVQSTK
NFINTFFENK'}

```

```

[38]: fastaSeqs = ">CCDS2.2|Hs109|chr1
MSKGILQVHPPICDPCGCRISSPVNRGRLADKRTVALPAARNLKKERTPSFSASDGSDG
SGPTCGRRPGLKQEDGPHIRIMKRRVHTHWDVNISFREASCSQDGNLPTLISSVHRSRHL
VMPEHQSRCEFQGRSLEIGLRPAGDLLGKRLGRSPRISSDCFSEKRARSESPQEALLPR
ELGSPMAPEDHYRRLVSALSEASTFEDPQRLYHLGLPSHGEDPPWHDPPHHLPSHDLLRV
RQEVAAAAALRGPSGLEAHLPSSTAGQRRKQGLAQHREGAAPAAAPSFSERELPPPPLLS
PQNAPHVALGPHLRPPFLGVPSALCQTPGYGFLPPAQAEFAWQQELLRKQNLARLELPA
DLLRQKELESARPQLLAPETALRPNDGAEELQRRGALLVLNHGAAPLLALPPQGGPGSGP
PTPSRDSARRAPRKGPGPASARPESEKEMTGARLWAQDGSEDEPPKSDGEDPETA AVG
CRGPTPGQAPAGGAGAEGKGLFPGSTLPLGFPYAVSPYFHTGAVGGLSMDGEEAPAPEDV
TKWTVDDVCSFVGGLSGCGEYTRVFREQGIDGETLPLLTEHLLTNMGLKLGPAKIRAQ
VARRLGRVFYVASFPVALPLQPPTLRAPERELGTGEQPLSPTTATSPYGGGHALAGQTSP
KQENGTLALLPGAPDPSQPLC
>CCDS3.1|Hs109|chr1
MAAAGSRKRRLAELTVDEFLASGFDSESESESESENSPQAETREAREAAARSPDKPGGSPSAS
RRKGRASEHKDQLSRLKDRDPEFYKFLQENDQSLNFSDDSDSSEEEGPFHSLPDVLEEA
SEEDGAEEGEDGDRVPRGLKGKKNVPVTVAMVERWKQAAKQRLTPKLFHEVVQAFRAA
VATTRGDQESAEANKFQVTDSAAFNALVTFCIRDLIGCLQKLLFGKVAKSSRMLQPSSS
PLWGKLRVDIKAYLGSAILVLSCLSETTVLAAVLRHISVLVPCFLTFPKQCRMLLKRMVI
VWSTGEESLRVLAFLVLSRVCRRHKDTFLGPVLKQMYITYVRNCKFTSPGALPFISFMQW
TLTELLALEPGVAYQHAFLYIRQLAIHLRNAMTTRKKETYQSVYNWQYVHCLFLWCRVLS
TAGPSEALQPLVYPLAQVIIGCIKIPTARFYPLRMHCIRALTLLSGSSGAFIPVLPFIL
EMFQQVDFNRKPGRMSSKPIFNSVILKLSNVNLQEAYRDGLVEQLYDLTLEYLHSQAHC
IGFPELVLPVVLQLKSFLRECKVANYCRQVQQLLGKVQENSAYICSRQRVSVFGVSEQQA
VEAWEKLTREEGTPLTLYYSHWRKLRDREIQLEISGKERLEDLNFPEIKRRKMADRKDED
RKQFKDLFDLNSSEEDDTGFSERGILRPLSTRHGVEDDEEDEEEDSSNSSEDGPDPA
EAGLAPGELQQLAQGPEDELEDLQLEDD
>CCDS4.1|Hs109|chr1

```

```

MGNSHCVPQAPRRLRASFSRKPSLKG NREDSARMSAGLPGPEAARSGDAAANKLFHYIPG
TDILDLENQRENLEQPF LSVFKKGRRRV PVRNLGKV VHYAKVQLRFQHSQDVSDCYLELF
PAHLYFQAHGSEGLTFQGLLPLTELSVCPLEGSR EHA FQITGPLPAPLLVLCPSRAELDR
WLYHLEKQTALLGGPRRCHSAPPQRR LTRLRTASGHEPGGSAVCASRVKLQHLPAQE QWD
RLLVLYPTSLAIFSEELDGLCFKGELPLRAVHINLEEKEKQIRSF LIEGPLINTIRVVCA
SYEDYGHWLLCLRAVTHREGAPPLPGAESFPGSQVMGSGRGSLSGGQTSWDSGCLAPPS
TRTSHSLPESSVPSTVGCSSQHTPDQANS DRASIGRRRTELRRSGSSRSPGSKARAEGRG
PVTPLHLDLTQLHRLSLESSPDAPDHTSETSHSPLYADPYTPPATSHRRVTDVRGLEEFL
SAMQSARGPTSSPLPSVPVSPASDPRSCSSGPAGPYLLSKKGALQSRAAQRHRGSAKD
GGPQPPDAPQLVSSAREGSPEPWLP L TDGRSPRRSRDPGYDHLWDETLSSSHQKCPQLGG
PEASGGLVQWI
>CCDS5.1|Hs109|chr1
MAADTPGKPSASPMAGAPASASRTPDKPRSA AEHRKSSKPVMEKR RRRARINESLAQLKTL
ILDALRKESRRHSKLEKADILEMTVRHLRSLRRVQVTAALSADPAVLGKYRAGFHECLAE
VNRFLAGCEGVPADVRSRLLGHLAACLRQLGPSRRPASLSPAAPAEAPAEVYAGRPLLP
SLGGPFPLAPLLPGLTRALPAAPRAGPQGPGGPWRPWL R
>CCDS6.1|Hs109|chr1
MGWDLTVKMLAGNEFQVSLSSMSVSELKAQITQKIGVHAFQQRLAVHPSGVALQDRVPL
ASQGLGPGSTVLLVVDKCDEPLSILVRN NKG RSSTYEVRLTQTVAHLKQQVSGLEGVQDD
LFWLTFEGKPLEDQLPLGEYGLKPLSTVFMNLRLRGGGTEPGGRS
>CCDS7.2|Hs109|chr1
MALRHLALLAGLLVGVASKSMEN TAQLPECCVDVVGVNASC PGASLCGPGCYRRWNADGS
ASCVR CGNGLTPAYNGSECRSFAGPGAPFPMNRSSGTPGRPHPGAPRVAASLFLGTFFIS
SGLILSVAGFFYLKRSSKLPRACYRRNKAPALQPG EAAAMIPPPQSSVRKPRYVRRERPL
DRATDPAAFPGEARISNV
>CCDS8.1|Hs109|chr1
MGSSQEGLRCQPSQPDHDADGHCGPDLEGAERASATPGPPG LLNSHRPADSDDTNAAGP
SAALLEGLLLGGGKPSPHSTRPGPFYIGGSNGATI ISSYCKSKGWQRIHDSRRDDYTLK
WCEVKSRDSYGSFREGEQLLYQLPNNKLLTTKIGLLSTLRGRARAMSKASKVPGGVQARL
EKDAAAPALEDLPTSPGYLRPQRVLRMEEFFPETYRLDLKHEREAFFTLFDETQI WICK
PTASNQKGKIFLLRNQEEVAALQAKTRSMEDDPIH HKTPFRGPQARVVQRYIQNPLLVDG
RKFDVRSYLLIACTTPYMIFFGHGYARL TSLYDPHSSDLGGHLTNQFMQKKSPLYMLLK
EHTVWSMEHLNRYISDTFWKARGLAKDWVFTTLKVRPLCPPVWE""

```

[39]: `protein_dict(fastaSeqs)`

[39]: {'CCDS2.2|Hs109|chr1': 'MSKGILQVHPPICDCPGCRISSPVNRGRLADKRTVALPAARNLKKERTPSFSASDG
DSDGSGPTCGRRPGLKQEDGPHIRIMKRRVHTHWDVNISFREASCSQDGNLPTLISSVHRSRHLVMPEHQSRCEFQRGSL
EIGLRPAGDLLGKRLGRSPRISSDCFSEKRARSESPQEALLPRELGPSMAPEDHYRRLVSALSEASTFEDPQRLYHLGL
PSHGEDPPWHDPHHLPSHDLRLVRQEVA AAAALRGPSGLEAHLPSSTAGQRRKQGLAQHREGAAPAAAPSFSERELPP
PLLSPQNAPHVALGPHLRPPFLGVPSALCQTPGYGFLPPAQAEFWAQQELLRKQNLARLELPADLLRQKELESARPQLL
APETALRPNDGAEELQRRGALLVLNHGAAPLLALPPQGPPGSGPPTPSRDSARRAPRKGGPGPASARPSESKEMTGARLW
AQDGSEDEPPKSDGEDPETA AVGCRGPTPGQAPAGGAGAEGKGLFPGSTLPLGFYAVSPYFHTGAVGGLSMDGEEAPA
PEDVTWKTVDDVCSFVGGLSGCGEYTRVFREQGIDGETLPLLTEEHLLTNMGLKLGPALKIRAQVARRLGRVFYVASFPV
ALPLQPPTLRAPERELGTGEQPLSPTTATSPYGGGHALAGQTS PKQENGTALLPGAPDPSQPLC',
'CCDS3.1|Hs109|chr1': 'MAAAGSRKRRLAELTVDEF LASGFDESESESESESPQAETREAREAAARSPDKPGGS
PSASRRKGRASEHKDQLSRLKDRDPEFYKFLQENDQSLNFSDDSSSEEEGPFHSLPDVLEEASEEEDGAEEGEDGDRV

PRGLKGKKNSVPVTVAMVERWKQAAKQRLTPKLFHEVVQAFRAAVATTRGDQESAEANKFQVTDAAFNALVTFCIRDLI
 GCLQKLLFGKVAKDSSRMLQPSSSPLWGKLRVDIKAYLGSAIQLVSCLETTVLA AVL RHISVLVPCFLT FPKQCRMLLK
 RMVIVWSTGEESLRVLAFLVLSRVCRHKKDTFLGPVLKQMYITYVRNCKFTSPGALPFISFMQWTLTELLALEPGVAYQH
 AFLYIRQLAIHLRNAMTTRKKETYQSVYNWQYVHCLFLWCRVLSTAGPSEALQPLVYPLAQVIIGCIKLIPTARFYPLRM
 HCIRALTLLSGSSGAFIPVLPFILEMFQQVDFNRKPGRMSSKPIINFSVILKLSNVNLQEKAYRDGLVEQLYDLTLEYLHS
 QAHCIGFPELVLPVVLQLKSFLRECKVANYCRQVQQLLGKVQENSAYICSRQRVSVFGVSEQQAVEAWEKLTREEGTPLT
 LYYSHWRKLRDREIQLEISGKERLEDLNFPEIKRRKMADRKDEDRKQFKDLFDLNSSEEDDTGEFSERGILRPLSTRHGV
 EDDEEDEEEGEEDSSNSEDGDPDAEAGLAPGELQQLAQGPEDLEDLQLSEDD' ,
 'CCDS4.1|Hs109|chr1': 'MGNSHCVPQAPRRLRASFSRKPSLKGNREDSARMSAGLPGPEAARS GDAAANKLFH
 YIPGTDILDLENQRENLEQPFLSVFKKGRRRVPVRNLGKV VHYAKVQLRFQHSQDVSDCYLEL FPAHLYFQAHGSEGLTF
 QGLLPLTELSVCPLEGSREHAFQITGPLPAPLLVLCPSRAELDRWLYHLEKQTALLGGPRRCHSAPPQRRLTRLTASGH
 EPGGSAVCASRVKLQHLPAQEQQWDRLLVLYPTSLAIFSEELDGLCFKGELPLRAVHINLEEKEKQIRSFLIEGPIINTIR
 VVCASYEDYGHWLLCLRAVTHREGAPPLPGAESFPGSQVMGSGRGLSSSGGQTSWDSGCLAPPSTRTSHSLPESSVPSTV
 GCSSQHTPDQANS DRASIGRRRTELRRSGSSRSPGSKARAEGRPVTPHLDLTQLHRLSLESSPDAPDHTSETSHSPLY
 ADPYTPPATSHRRVTDVRGLEEFLSAMQSARGPTPSSPLPSVPVSPASDPRSCSSGPAGPYLLSKKGALQSRAAQRHRG
 SAKDGGPQPPDAPQLVSSAREGSPEPWLPLTDGRSPRRSRDPGYDHLWDETLSSSHQKCPQLGGPEASGGLVQWI' ,
 'CCDS5.1|Hs109|chr1': 'MAADTPGKPSASPMAGAPASASRTPDKPRSAAEHRKSSKPVMEKRRRRARINESLAQ
 LKTLILDALRKESSRHSKLEKADILEMTVRHLRSLRRVQVTAALSADPAVLGKYRAGFHECLA EVNRFLAGCEGVPADVR
 SRLLGHLAACLRQLGPSRRPASLSPAAPAEAPAEVYAGRPLLP SLGGPFLLAPLLPGLTRALPAAPRAGPQGPGGPW
 RPWLR' ,
 'CCDS6.1|Hs109|chr1': 'MGWDLTVKMLAGNEFQVSLSSMSVSELKAQITQKIGVHAFQQLAVHPSGVALQD
 RVPLASQGLGPGSTVLLVVDKCEPLSILVRNNKGRSSTYEVRLTQTVAHLKQVSGLEGVQDDLFWLTFEGKPLEDQLP
 LGEYGLKPLSTVFMNLR LRGGGTEPGGRS' ,
 'CCDS7.2|Hs109|chr1': 'MALRHLALLAGLLVGVASKSMEN TAQLPECCVDVVGVNASC PGASLCGPGCYRRWN
 ADGSASCVRGNGTLPAYNGSECRSFAGPGAPFPMNRSSGTPGRPHGAPRVAASLFLGTFFISSGLILSVAGFFYLKRS
 SKLPRACYRRNKAPALQPGEAAAMIPPPQSSVRKPRYVRERPLDRATDPAAFPGEARISNV' ,
 'CCDS8.1|Hs109|chr1': 'MGSSQEEGLRCQPSQPDHDADGHCGPDLEGAERASATPGPPGLLN SHRPADSDDTN
 AAGPSAALLEGLLLGGGKPSPHSTRPGPFYIGGSNGATI ISSYCKSKGWQRIHDSRRDDYT LKWCEVKS RDSYGSFREG
 EQLLYQLPNNKLLTTKIGLLSTLRGRARAMSKASKVPGGVQARLEKDAAAPALEDLPTSPGYLRPQRVLRMEEFFPETY
 RLDLKH EREAFFTLFDETQIWICKPTASNQKGIFLLRNQEEVAALQAKTRSMEDDPIHHKTPFRGPQARVVQRYIQNPL
 LVDGRKFDVRSYLLIACCTPYMIFFGHYARLTSLYDPHSSDLGGHLTNQFMQKKSPLYMLLKEHTVWSMEHLNRYISD
 TFWKARGLAKDWVFTTLKVRPLCPPVWE' }

[]: