

Задача для курсовой работы - написать и отладить программу из двух частей: А) создающую и обновляющую простые базы данных по нуклеотидным последовательностям интронов из заданных видов растений; Б) исследующую заданную нуклеотидную последовательность на предмет возможности вставить в нее интрон(ы) с имеющимся нуклеотидным контекстом из созданных баз.
Необходимые знания: программирование (желательно, C++), биоинформатика.

К кому обратиться:

1. Алексей Фолимонов

Старший научный сотрудник

Группа молекулярной биологии вирусов растений

Центр "Биоинженерия" РАН

Проспект 60-летия Октября, д.7, к.1

Москва, 117312

телефон +(8) 499-135-12-40

мобильный +(8) 909-672-73-64

e-mail folimonov@biengi.ac.ru

2. Алексей Анатольевич Аграновский

Профессор кафедры вирусологии Биологического факультета МГУ

телефон +(8) 495-939-23-63

e-mail aaa@genebee.msu.su

Вот, что делаю я (а должна делать программа А):

1. Найдем в ГенБанке все последовательности *Arabidopsis thaliana*, содержащие интроны:

The screenshot shows a web browser window displaying the NCBI Entrez Nucleotide search results for the query "intron AND Arabidopsis thaliana". The search results are sorted by file name, and the first 8 results are listed. The first result is AY578790, which is SEN2 (SPlicing ENdonuclease 2). The page also includes a "Recent Activity" sidebar and a taskbar at the bottom.

Accession	Report	Organism	Links
AY578790	Reports	Arabidopsis thaliana [gi:46397600]	Links
X89366	Reports	A. thaliana DNA fo... [gi:897676]	Links
BN000023	Reports	TPA_exp: Arabidop... [gi:28317384]	Links
AK317057	Reports	Arabidopsis thali... [gi:222423560]	Links
AK316802	Reports	Arabidopsis thali... [gi:222423073]	Links
M80567	Reports	A. thaliana non-sp... [gi:1174373]	Links
L24070	Reports	Arabidopsis thali... [gi:456015]	Links
T22568	Reports	Arabidopsis thali... [gi:4046671]	Links

Или, скажем, последовательности *Nicotiana tabacum*:

The screenshot shows the NCBI Nucleotide search interface. The search query is 'intron AND "Nicotiana tabacum"[porgn: txid4097]'. The results page displays a list of 127 items, with the first 9 items visible. The items are listed with their accession numbers, reports, and links. The first item is X14059, and the last is FJ664246. The page also includes a 'Recent Activity' sidebar and a 'My NCBI' section.

Или любого другого растения, для которого есть название и txid, хотя арабидопсис, конечно, интереснее, просто потому, что последовательностей больше...

2. Сохраним номера последовательностей в соответствующих файлах.
3. Посмотрим на последовательность M84466 из «табачного» файла:

The screenshot shows the NCBI Nucleotide page for the Tobacco phenylalanine ammonialyase (tpa1) gene. The page displays the gene's accession number (M84466.1) and its complete coding sequence (cds). The page includes a 'Features' section with details about the gene's location, definition, and source. The 'Sequence' section shows the nucleotide sequence. The page also includes a 'Recent Activity' sidebar and a 'My NCBI' section.

Интересует вот эта информация (выделено синим):

The screenshot shows the NCBI Nucleotide database entry for the tobacco phenylalanine ammonia-lyase (tpal) gene. The page displays genomic DNA coordinates, CDS coordinates, and exon/intron structures. A specific intron sequence (2152..4083) is highlighted in blue. The browser interface includes a search bar, navigation buttons, and a taskbar at the bottom.

4. Выделяем последовательности интрона (показано красным) и по 20 нуклеотидов 5'- и 3'-концевого экзона (показано зеленым), вот они:

```

cttcaaaagaacttattag gtaaacccac
2161 ttttcatcta ctttaactttg ctttctgcac cagaatttta atattttggt tctctctaaa
2221 gattttttatc ctttaactttg ttaaataatt taaattaatt attttcaaga tatcagagt
2281 tattatacac aagattttctg tcagaagtgt gcagtgacc ttagtcgaa aaaaaata
2341 gcatctacat ttttctatga cagcgttaag gaaattttac acggtcagta gaaaagggtta
2401 ttttaagttta ttattttagt ttattgtaag cagctatata tatatata tataatata
2461 gttatgttt agatggtcgt aaatttatt tgtgtgacc atgtatagta atagttaaat
2521 tcttattctg caagaaaatc tgactagcat ttatgagtt atgactact attttgagtt
2581 aataactatt aatgatata actaatttga cttataaaa aatatctta ttatttgag
2641 tcaatagtt ttaatgatt actttaactc ttatgtcaac tactctttt gtaagaggaa
2701 atgagttta aatgatata actttaaggt aaaatgtgta cattttctat gtaatctaaa
2761 cgattatagc agctatatt caagtatta ttccgacca aatatgaaa gctcttatt
2821 atatagaac aatataatct ggtactata atttttctta tttatcata gtagagaga
2881 taaagtttgt atatgtaaac atttatcaca cgcgctagac caagaactaa gtaacttga
2941 tttacatctt gaattttgaa ttattgagta aagaaaaata ttaaatgccc ctggaaaaa
3001 tttttatata tggaaacttga tacgaagtac aggaattttg catctacatc tataaatcta
3061 attatataaa taaattttg tacaaactgta aaaaaaac aagccactag ttgactttt
3121 gttgggttg aaaaagagt ggtgggttgg ggaggagatt ttagtcattt caacttctt
3181 ctagttttgt ataaagtaca tgaagaatt tattcagtaa atttttcta ctactattt
3241 cctaaattta gtaagaaaaa ataatttaa taacatctt taaatagtt aattttttt
3301 gttgggtctc aaataataga tgagagacaa gtgggcttat taaacagtt taaattcca
3361 actatttggg tcgaatccc gtagtcccat gtgctgtctc ttgtcagata tctatattt
3421 taaagtaca aataattttt aaaacaaata ttgcataaat aagtaatttt atactattt
3481 aaagagaatt gtaaaaaaa ttatgaaatg acaatttta atctcaaat taaacttca
3541 tttattttt gggcctaatt taccttaac ctgcccga actattata tctgtagct
3601 acaaaagtgt aaaaacttta tatattttt tataacaac acatataat atatacaaaa
3661 atatatatg attttgcgtt attattttt tagcggctat gtcattttt tctttttt
3721 attttccct cccgtgcca caaggccaa aatatacaa atgtgtgcc cactttaaag
3781 gggcttata tcaactgaaa tgaaggagt tgatgtaaca cactacaagt ttttctttt
3841 ttttctaa ccttataat ccttgatccc tattctagtc attttctgt actaacaat
3901 atatatcca atataatct atactaatt tctcagac cctcggctaa ttttctata
3961 ctagtaatt taaatatcg tcaaaaagta gataggcaa agtacaata aacaagtgc
4021 ttatttaaat ttggttgga caattcagtt caatagctg tgtgtgtgac attgaattg
4081 cag gttcttgaactgctgtggtt

```

5. В базе данных нам нужна информация: 1) о последовательности экзонов; 2) о последовательности и длине интрона; 3) все, что можно получить из «ACCESSION M84466», откуда мы сей интрон почерпнули, а там можно найти ссылки, сорт растения, экотип и т.д.; 4) природный кодоновый (аминокислотный) контекст:

Экзон 5' Экзон 3'
ААСТТАТТАГ | ГТТСТТГААТ
N L L G S *
T Y * V L E
L I R F L N

Например, в этом случае в природной мРНК – контекст **LI {R} FLN**, т.е. кодируемый триплет аргинина **AGG** «разрывается» **AG | G**. Природный контекст важнее, поскольку в «искусственном» (два других могли бы иметь право на существование) есть потенциальная проблема с «codon pair bias».

Это все о программе А.

Программа Б:

1. Ограничим кодоновый контекст «табачной» (любой) базы, кодируемыми аминокислотами, возьмем только по одной аминокислоте (это число может меняться) «справа» и «слева» от центральной: **I {R} F**.
2. Найдем в заданной аминокислотной последовательности (например, в репликазе вируса желтухи свеклы) эти три аминокислоты:

MAFLNVS AVPS CAFAPAFAPHAGAS PIVPD SFPCV PRY SDD ISHFRLT LSLDF SVPRPLS LNARVHLSAS TDNPLPS LPLGFHAE TFVLE LMGSSAP
FSI PSRHI DFVNRPF SVFPT EVL (1) SVS SLRTP SRL FALLCDF FLYCSKPG PCVEIAS FST PP PCLVSNCAQ IP THAEMES IRE (1) PTKLPA
GRFLQFHKKRY TKRPE TL I IHESGLALKTS ALGVT SKPNSRPI TVKSASGEKYEAEYI SRKDFERSRRRQQT PRVRS HKPRK INKA VE PFFF PEE PK
KDKRRAS LPT EDEGE IT FTGLRFP LSETPKEE PRLPKFREVE IPVVKHHA VPAVSK PVRTFRP VAT TGAE YVNARNQC SRRPRNH PILRSASYTF
GFKMPLQREMFKEKKEYYVRS (1) KVVSS CSVTKSPLEALAS ILKNLPQYSYNSERLKFYDHF I GDD FE IEVHP LRGKLSVLL IL PKGEAYCVVT
AAT PQY HAALT IARGDRP RVGELQYRPGEGLCYLAHAALCCALQKRT FREED FFVGMYP TKFVF AKRLTEK LGP SALKH PVRGRQV SRS LFHCDVA
SAFSSP FYSLPRF IGGVEEAEPE IT SSLKHKAI ES VYERV SIHKD NLL ARSVEKD LID FKDEI KSLSK EKRS (2) VTVPFYMGEAVQ SGLTRAYPQF
NLS FTHSVYSDHPAAAGSR LLENET LASMAKSS FSDIGCC PLFHI KRGSTDYHVC RPI YDMKDAQRRV SRELQARGL (1) VENLSREQLVEA QARVS
VCPHTL GNCNVKSDVL IMVQVYDAS I NE IA SAMVLKES KVA (1) YLTMVT PGE LLDEREA FAI DALGCDVVVDTRRDMVQYKFGSSCYCHKL SNIKS
IMLTPAFT FSGNLF SVEMYENRMGVNYYKI TRSAY SPE IRGVKTLRYRRACTE VVQVQLPREDKT LKT FL SGYDY IY LDKAFVSRVFDYVVS NCSV
NSKTFEWWVSY IKSSKSRVVI SGKV IHRDVHIDLKHSECF AAVMLAVGVR SRTTTEFLAKNLYY TGDASC FET IRE (2) LFREWSRRAYAE INRSF
RKLMS ILSAGLD YEF LD LDNSLQHLL EYS EVEVRVSI AQNGEVD CNEENRVL TE IIAEAADRKS IAQGLSGALS SVPTQ PRGGLRGSSRS GV (1)
SFLYLVEE VGNL FFSVGDVAVRFLVKVFKT FSDSP I FRVVRMFLDLAEASPFFVSVVSLCAWLREAVSAFSSWVADR TVSES VKTFVNRVTKRFLNF
MSAKTLTKKFFRF FLSASALAKTVVRKAKV ILEAYWEVWFESI LSDSGEYSAVEFCSSV I TLLTNSGRLLPGFS PAIIT EVL (2) LDLATKIS IE
VL (3) LKQISPADSTASSALYRRVLSEI LSNFR TMGEHGI FTKVFLLCGF LPV FVRKVA LCVPGDMATYARFLE YGVDD LFFLGRSVNS IKNYLCV
VAAGLVDS IVD SVLKLKLSGV (2) AKERVLGFKSKI IKN FLNVFRKAKVVT RTS SSTDLSEDEYFS CDE SKPGLRGSS SRFTL SRL LD IFFNF LKSK
LVI ENACFSAYER IERNMKLYFFPLNSSEE EARRL IRCAGDFD YLSDSAFDEDEMLRQAF EQYYS SDD ESVT YDGKPTVLRSLY NVSRRFLE TFCNG
PKF FVKVSNYFKALYSRLLRVLPWDRNLS DSPGLKGGNEKALAKFFKT CVI TACECVS QIC CLRLI RL CWGTPACGLVRLFYI TYSTRVLSRVV
VAVAVC PLLVRNE LDGLSDGL TNMGVSVFRRLFVALRRAL SAY SNSALRRKI IEF IFGNI HHPFDVAV IE TNEVAPE PLSPEVDI DVDCDFGSDSES
VSSDEVASNPR PGLHGGSSRS SNFL TSLVKVVKLARR IPRLLFR LRNFVAYFVE RRLAS KRLKT FIGLARLFDNFSLTS VVYLLQE YDSVLNAF ID
VELILLNSGSVNVLPVSVWRGSLTKLAEAVGSGFAS FLGRMCCRVSDWCSS SSNAGCNFMS PVRTKGFV PPS SSGSTASMYERLEALES DIREH
VLS TCRVGSDEEEERPKEVTE PGI EHTSEDVVP IRSHS QPLSGGECYS SEDRE ENERANL LPHVSKIVSERRGL (2) ETARRNKRTLHGVSEFLNAI
NTSNEQPRPI IVDHSPESRAL TNSVREFYY LQELALFE LSKLRE YYDQL KVA (2) NFNROECLCDKDEDMFVLRAGQGVVS GRNSR LPLKHFKGHE
FCFRSGGLVPYDGT SRVD TIFHTQT NFVSANALLS GYLSYRTFT TNL SANVLLYEAP PGGKTT TLI KVCETF SKVNS LI LTANKSSREE ILAKV
NRI VLDEGDTPLQTRDRI LTI DSYLMNNGRL (3) TCKVLYLDECFMVHAGA AVAC IEF TKCDSA I LFGDSRQ IRYGRCELD TAVLS DLNRFVDES
RVYGEVSYRCPWDVCAWLS TFYPKTVAT TNLVSAQSSMQVRE IE SVDDVEYS SE FYLLTMLQSEKDLLKSGKRS (3) RS SVEKPTVL TVHEA QG
E (1) TYRKVNLVR TKFQEDDF FRSENI TVALS RHVES LTYSVLS SKRDDATAQA IVKAKQLVDA YRVYPT SFGGSTLDVSVNPS TSDRSK KAS SA
PYEVINSFLESVVP GTTSVDFGDVSEEMGT QVFES GADNVVIRDS A PVNKST'DHD PQRV*

Есть такие! Встречаются дважды (показано красным на желтом фоне)! Остальные «подчеркивания» - «интроны» арабидопсиса.

3. Для сканирования заданной аминокислотной последовательности на предмет наличия «аминокислотного» контекста интрона из заданной базы (или баз), неплохо задавать предельный размер потенциального интрона.

Так работаю я, а должна работать программа Б.

4. А и Б могут быть совмещены в одном интерфейсе.