I have installed Ubuntu (Linux OS) on windows computer as a dual boot --- both windows OS and Ubuntu OS can be loaded at start. Such installations of Linux are much more powerful for computation over virtual machine.

Please see guide "2.1_Install_Ubuntu_dual_boot_any_computer.pdf" and "2.2_Install_Ubuntu_different_software.pdf" for detailed installation steps (with screenshots and/or video descriptions where available).

**end of SM note 1**

**SM note 2:**

What IDE (integrated development environment) or Editor do I need? How to run/execute Perl script?

Answer: After reviewing multiple options,

IDE point of view Padre and Epic it seems both are good choices for Execution + Syntax and also Edit (ref: https://www.dunebook.com/best-perl-ide-and-editors/). Disadvantage of IDE unlike Editors is that a) IDE are heavy weight applications b) Don't support multiple languages so got to learn again each time.

Editor point of view that can do not only Edit but also does Execution + Syntax support, github's Atom is great cause a) Atom is light weight application unlike IDE b) Atom is more amenable to multiple languages and features so there's the comfort level that comes from using Atom previously.

- For Syntax on Atom: install "language-<program language>" package "language-perl" from withing Atom using >Packages > Settings View > Install Packages/Themes > search package name > click install. This "adds syntax highlighting and snippets to <program language> files in Atom". The file extension tells atom what language it is, .pl for Perl, .py for Python, .js for JavaScript and .java for Java.
  Perl "language-perl" https://atom.io/packages/language-perl, Python "language-python" https://atom.io/packages/language-python, JavaScript "language-javascript" https://atom.io/packages/language-javascript, Java "language-java" https://atom.io/packages/language-java and other languages too.
- For Execution on Atom: Atom can also execute/run scripts if we install package "script" from withing Atom, to run script file click Packages > Script > Run Script or shortcut Ctrl+Shift+B, output shows in tiny space below script file, if click 'show in new tab' icon it shows the output in new tab of Atom.
  Perl https://youtu.be/KzEuTjNwzvk, Python https://youtu.be/QuiquBXdA1o, JavaScript https://youtu.be/Aj7Iah4ksH8, Java https://youtu.be/Wg5Peun14YM and other languages too.
- For Terminal window on Atom: On Atom we can type in terminal window install package "terminal-plus" from withing Atom https://atom.io/packages/terminal-plus.
- Additional IDE features for Atom: Its not essential I think, but additional packages are available on Atom for more IDE type features. To find this type "Atom <language name> IDE" in search engine and follow instructions.

Please see guide "2.3_Install_Atom_Its_More_Than_Just_Text_Editor.pdf" for detailed installation steps (with screenshots and/or video descriptions where available).

**end of SM note 2**

# Background

We provide gene annotation for the human genome.
The latest assembly available is GRCh38 but some users still use the previous assembly, GRCh37, and want to convert coordinates between the two assemblies.
You can find more information about genome assemblies here:
https://www.genome.gov/12011238/an-overview-of-the-human-genome-project/

For this exercise, we use the latest Ensembl release, 104.

**SM note 3:** This question is about conversion of gene coordinates between two assemblies.
https://www.ncbi.nlm.nih.gov/assembly/88331

✅ GRCh38 - hg38 - Genome - Assembly - NCBI

https://www.ncbi.nlm.nih.gov/assembly/88331
1. GRCh38 Genome Reference Consortium Human Build 38 Organism: Homo sapiens (human) Submitter: Genome Reference Consortium Date: 2013/12/17 Assembly type: haploid-with-alt-loci Assembly level: Chromosome Genome representation: full Synonyms: hg38 GenBank assembly accession: GCA_000001405.15 (replaced) RefSeq assembly accession: GCF_000001405.26 (replaced) IDs: 88331[UID] 883148 [GenBank ...

https://www.ncbi.nlm.nih.gov/assembly/2758/

✅ GRCh37 - hg19 - Genome - Assembly - NCBI

https://www.ncbi.nlm.nih.gov/assembly/2758
1. GRCh37 Genome Reference Consortium Human Build 37 (GRCh37) Organism: Homo sapiens (human) Submitter: Genome Reference Consortium Date: 2009/02/27 Assembly type: haploid-with-alt-loci Assembly level: Chromosome Genome representation: full Synonyms: hg19 GenBank assembly accession: GCA_000001405.1 (replaced) RefSeq assembly accession: GCF_000001405.13 (replaced) IDs: 2758[UID] 2468 [GenBank ...

**end of SM note 3**


# Ensembl Perl API

The Ensembl Perl API is used to generate annotation as well as access it programmatically.
You can find more information about it here:
http://www.ensembl.org/info/docs/api/core/core_tutorial.html

**SM note 4:** Setup and installation for Exercise. Ensembl Perl API tutorial (with video) http://uswest.ensembl.org/info/docs/api/api_installation.html was very helpful.
Presently, I did not make the additional installations for 'Variation genotype and frequency data'.

I used Ubuntu Release 20.04. Typing command 'lsb_release -a' on Terminal displays this information. *Please, note that a) text after '#' is for human not computer, b) text I greyed out is the command that needs to be typed on or copied to Terminal or Atom as the case maybe and c) after typing or copying the command hit Return/Enter key for it to execute.*
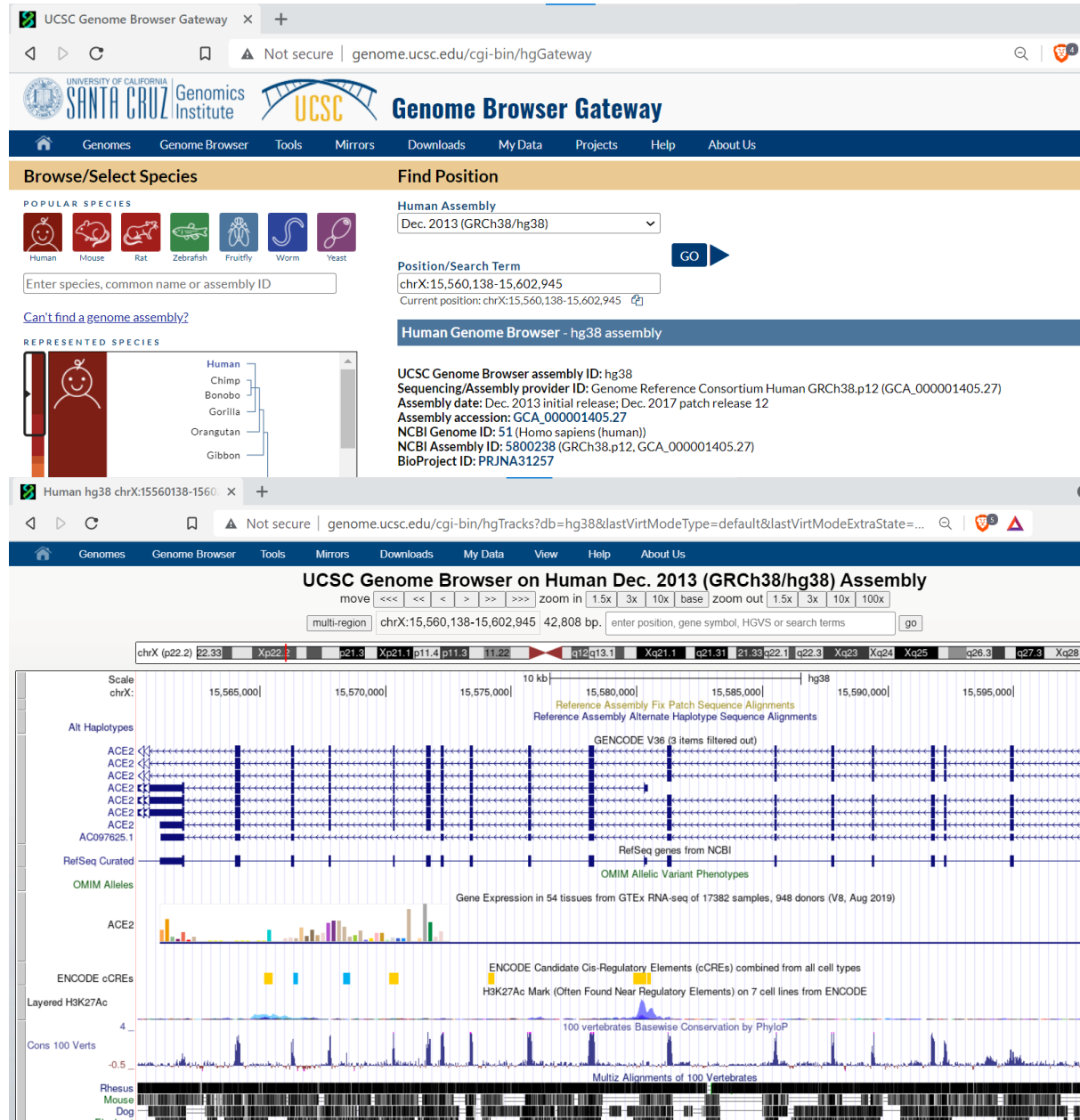lsb_release -a #Displays version of Ubuntu d) Clicking on Terminal window and pressing Ctrl+C or Cmd+C ends any code running on Terminal.

Please see guide "2.4_Install_Perl_PerlModules_BioPerl_EnsemblPerlAPI.pdf" for detailed Perl, Perl modules, BioPerl and Ensembl Perl API installation steps (with screenshots and/or video descriptions where available).

**end of SM note 4**

**SM note 5**

Used tutorials http://uswest.ensembl.org/info/docs/api/core/core_tutorial.html and http://uswest.ensembl.org/info/docs/api/general_instructions.html to get an idea of lines of code that need to be written to access and use Ensembl Perl API. I used Atom to write the 'convert_chr_region_hu.pl' script and executed the script from Ubuntu Terminal above where my .bashrc has PERL5LIB in its path. At this time, I don't know how to get Atom to run 'convert_chr_region_hu.pl' script using the installation libraries PERL5LIB. Using UCSC genome browser can view the region of the genome format "chr1:213941196-213942363"



Can also obtain sequence of DNA for the region by clicking > View > DNA

**end of SM note 5**

# Exercise

Using the Perl API on the latest human data for Ensembl release 104, convert coordinates on chromosome 10, from 25.000 to 30.000 to the same region in GRCh37.

**Git Link: https://github.com/smukher2/EMBL_git_ensembl_perl_API_Ex1.git**

**Point 1 of 4: Changed Input or Query Coordinates As Original One Provided In The Question Has No Sequence "No Data" or Gap region:**
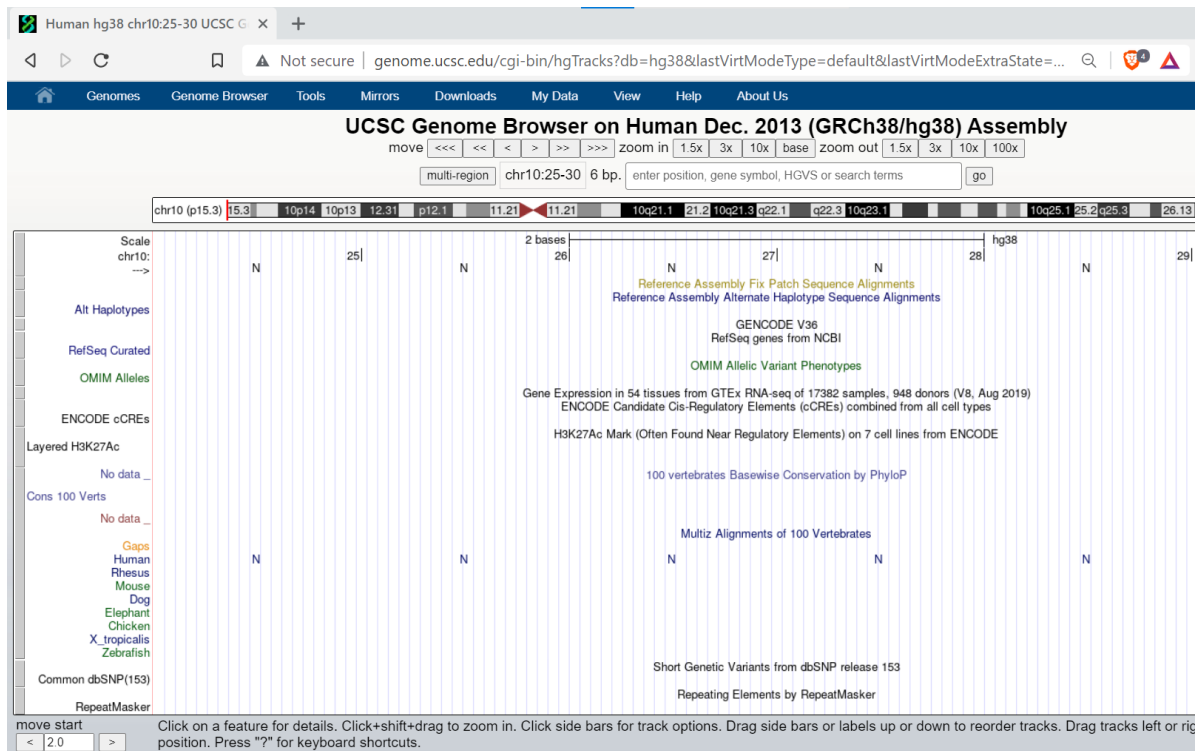
Upon checking UCSC genome browser for region GRCh38 chromosome 10, from 25.000 to 30.000, I found that this region does not exist.



```
>hg38_dna range=chr10:25-30 5'pad=0 3'pad=0 strand=+ repeatMasking=none
NNNNNN
```

Therefore, for conversion to GRCh37, I used most likely region GRCh38 chromosome 10, from 25000 to 30000 to continue with this exercise.



```
>hg38_dna range=chr10:25000-30000 5'pad=0 3'pad=0 strand=+ repeatMasking=none
TAAAAAATAATTTAAAGAAAATAGAGCTTGATAAATCTGTTTCTTCATCT
TCATAATAATGTCAAATTTGTTGAGATTTTTTTTAAATGGCACGATTTGT
CTACAGATCTTTGTACTCTGGCTTAAATTAAATGTATATACAACTTATAT
AATAAAATACTGGGCTGTTTTATTATTATATTGACAATTTCTTAGTCACC
ACTACTTTGATCATATTCTATAAATGGCACTGTGAGACCTATCACTGGTT
ATCAACAACATAAAGTGTTTATGACAATTTCATGTGACAGATGAAGGAAG
TTTGGTTAGTGCAGTGGCAATAGTTGGCAACCGAGAAACTGAGCTCACAA
TTTTAGAAAAGTATTATTCTTTCTAACTAGATATTTCCATGAAAAGAACC
TTGTGAAATGCAGAAATGCCAAAAGAATAATTATAAGATAATAACCAAAT
AACTAAAATACATCTCAATTTCTCTGACAATCTAACTTAATGCAGGGTAC
AGCCAGAATCATGATGTTTACAAACATAGACATTTACGATCTCGAGACAT
TTAATTATTAGCTGGTTGAGGAAAACAGTAGTAGAACTGTTATTGCAGAG
ATAAAGACAGAAAACAGTGGAAGGGCTTGATCACAAGTCATATGTTCCA
GTTTTGGAAACCTCTAAAAATCCTGCCTGGTGATAGGTTTTTTCTTATCTC
ATTATTGAAGAATATATGAGCAGGCAGGGCACGGTGGCTCAAACCTGTAA
```

4

Please see below screenshots of these genomic regions.
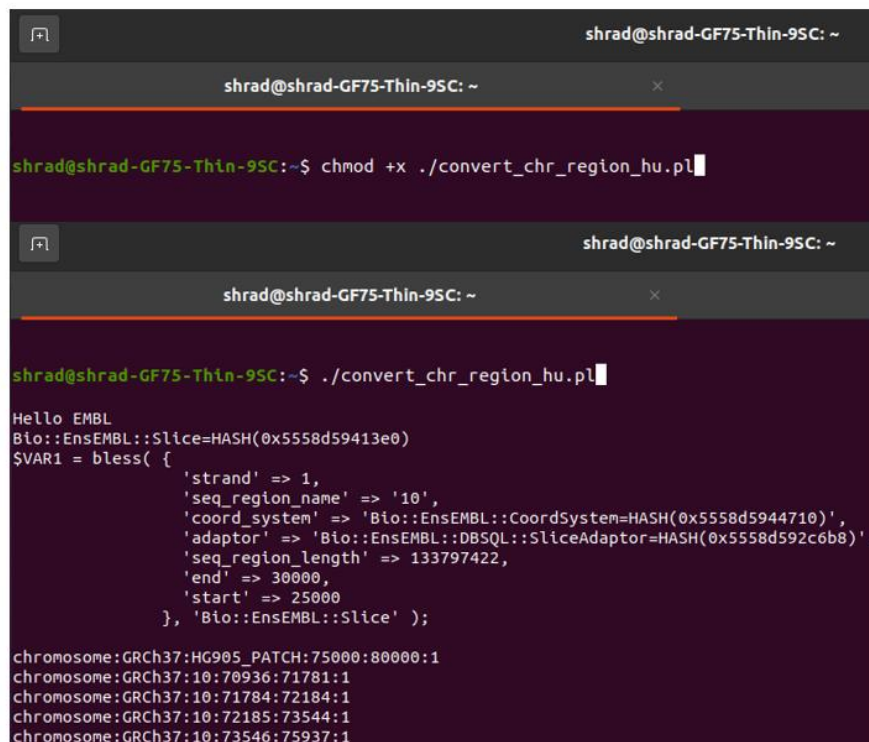GRCh38 chromosome 10, from 25.000 to 30.000 does not exist.



GRCh38 chromosome 10, from 25000 to 30000 does exist, and was used for this exercise.

**Point 2 of 4: Conversion of Coordinates form GRCh38 to GRCh37:** I used Ensemble Perl API to convert GRCh38 human genome region chr10: 25000-30000 to corresponding GRCh37. I used 'Atom' Editor to write the 'convert_chr_region_hu.pl' script and executed the script from Ubuntu Terminal using command (Ref: file permission https://cets.seas.upenn.edu/answers/chmod.html and running perl script https://www.geeksforgeeks.org/hello-world-program-in-perl/). I found helpful tutorials at https://m.ensembl.org/info/website/tutorials/grch37.html and https://rest.ensembl.org/documentation/info/assembly_map.

which perl #this gives location of perl used in first line of script for me its #!/usr/bin/perl

chmod +x ./convert_chr_region_hu.pl #gives file permission to make script executable for all users

./convert_chr_region_hu.pl #run the script



**Point 3 of 4: Results From Exercise**: The results for this exercise, for conversion to GRCh37 is as follows,
chromosome:GRCh37:HG905_PATCH:75000:80000:1
chromosome:GRCh37:10:70936:71781:1
chromosome:GRCh37:10:71783:72184:1
chromosome:GRCh37:10:72184:73544:1
chromosome:GRCh37:10:73545:75937:1

**Point 4 of 4: Pre-requisites for this exercise:** I started with the assumption that anyone trying to run this exercise has internet connection and a computer (Mac, Windows or Linux), with no background in coding. Therefore, I uploaded detailed description of all setup/installation steps and a detailed tutorial in github repository **Git Link: https://github.com/smukher2/EMBL_git_ensembl_perl_API_Ex1.git.**

Installation instruction files are:

2.1_Install_Ubuntu_dual_boot_any_computer.pdf

2.2_Install_Ubuntu_different_software.pdf

2.3_Install_Atom_Its_More_Than_Just_Text_Editor.pdf

2.4_Install_Perl_PerlModules_BioPerl_EnsemblPerlAPI.pdf

Tutorial file is:

4.1_Tutorial_Complete_Guide_Ubuntu_To_EnsemblPerlAPI.pdf


**end of SM Answer Exercise**


# Alternatives

Describe at least one other way of retrieving the same information, along with its advantages and disadvantages.


**SM Answer Alternatives**

**Git Link: https://github.com/smukher2/EMBL_git_ensembl_perl_API_Ex1.git**


**Point 1 of 5: Alternative Tool**: As an alternative we can use a 'LiftOver' Ensembl tool http://uswest.ensembl.org/info/docs/tools/index.html that is web-browser based and does not require installation http://uswest.ensembl.org/Homo_sapiens/Tools/AssemblyConverter.


**Point 2 of 5: Advantages Of Alternative Tool**: The Ensembl "LiftOver" tool has the following advantaged:

1. Web-based ready to use tool that does not require setup/installation.
2. Usage of tool does not require any coding experience.


**Point 3 of 5: Disadvantage Of Alternative Tool**: The Ensembl "LiftOver" tool has the following disadvantaged:

1. Web results are deleted after 10 days. So user needs to manually download results in 10 days. In the coding exercise results are automatically saved/downloaded to the computer.
2. Its tedious to use when user has multiple coordinate files to convert, such as from samples of different experimental and control groups, and replicates. For repetitive tasks coding exercise method is preferable.


**Point 4 of 5: Usage Of Alternative Tool**: The steps to do conversion of GRCh38 chromosome 10, from 25000 to 30000 to GRCh37.

1. Visit 'LiftOver' tool http://uswest.ensembl.org/Homo_sapiens/Tools/AssemblyConverter
2. Input coordinates GRCh38 in bed format and click 'Run' to execute conversion, as shown in screenshot https://m.ensembl.org/info/website/upload/bed.html.
   chr10  25000  30000

3. Click 'Download' to download GRCh37 converted results as shown in screenshot below.



**Point 5 of 5: Results From Alternative Tool**: The results for this exercise, for conversion to GRCh37 is as follows,

HG905_PATCH     75000  80000
chr10  70936  71781
chr10  71783  72184
chr10  72184  73544
chr10  73545  75937

end of SM Answer Alternatives