

Running the Del Sequencing Tool

This notebook serves as a replacement for the earlier sequencing data processing pipeline that was used previously for NADEL. Several improvements have been made, namely:

- An external web-based preprocessing step is no longer needed and has been directly integrated into the pipeline.
- The final output is a DataWarrior-ready csv file featuring either absolute or normalized sequence counts. No further manipulation in Excel is needed.
- The prior sequence searching algorithm did not attempt to correct error sequences. Some error sequences are able to be recovered using this processing method.
- The code in this workflow is ~500 lines shorter than the previous program and is annotated more clearly. It is also designed to be more modular and making modifications and updates to suit any library is considerably easier to manage.

Introduction –

This Del Screening Tool exists as a Jupyter Notebook and features several code cells that have been logically separated according to function. These cells are executed sequentially to move data through the full processing pipeline.

What does it do?

The tool takes a .fastq file and ultimately returns two .csv files containing normalized and absolute sequence counts organized according to BBs and screening conditions.

What Libraries are supported?

Currently, libraries with up to 4 encoding regions are supported. This translates to a maximum of 2 BB codes and 2 PCR codes. The code can be easily modified to support more codes should the need arise in the future.

Package Contents –

Name	Date modified	Type	Size
classes	3/13/2024 2:28 AM	File folder	
docs	3/4/2024 12:08 PM	File folder	
library	3/12/2024 10:18 PM	File folder	
outputs	3/12/2024 10:19 PM	File folder	
utils	2/28/2024 2:40 PM	File folder	
DelScreenProcessing.ipynb	3/14/2024 6:40 PM	Jupyter Source File	12 KB
ExampleData.fastq	3/4/2024 12:53 PM	FASTQ File	71,173 KB

- **DelScreenProcessing.ipynb** – The main notebook file used for data processing.
- **classes** – A subdirectory used as a python package for holding class modules.
- **library** – A subdirectory used as a python package for loading library-specific information.
 - **primers_all.py** – A file for storing a full list of PCR encoding sequences and IDs.

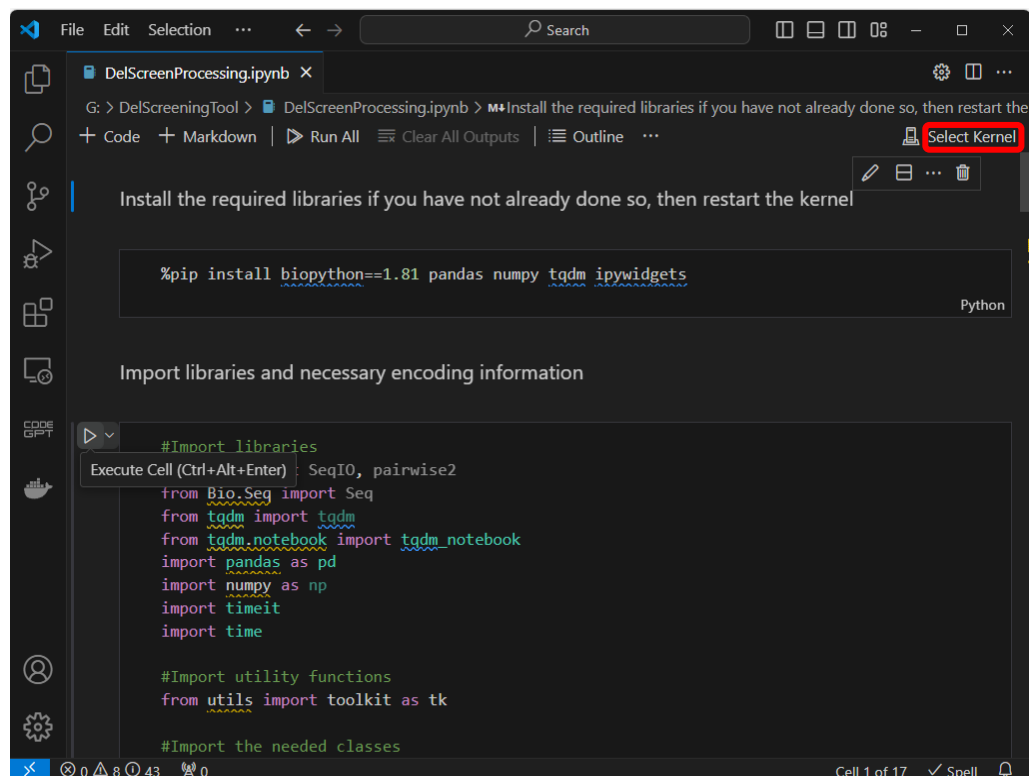
- **primers_screen.py** – A file for storing only the PCR encoding sequences and IDs used in a specific set of screens.
 - **fragment_codes.py** – A file for storing the encoding sequences and id's for the building blocks in each diversity element.
 - **fragment_smiles.py** – A file for storing the smiles and id's for the building blocks in each diversity element.
- **utils** – A subdirectory for additional utilities.
 - **DictFormatter.ipynb** – A notebook for auto formatting entries that can be pasted into the library .py files. Read more about this in the customization documentation.
 - **SeqView.ipynb** – A notebook for examining preprocessed fastq files. Useful for determining the slice indices for the encoding sequences.
 - **toolkit.py** – A file containing a variety of utility functions used by the main notebook.
 - **ExampleCodes.csv** – A csv file used as an example input for the DictFormatter notebook.
 - **ExampleData_Preprocessed.txt** – A txt file used as an example input for the SeqView notebook.
- **docs** – Documentation for the setup and usage of the included files.
- **outputs** – Files generated by the main notebook are organized and output here.
- **ExampleData.fastq** – A small .fastq file containing 200,000 NADEL reads that can be used for testing.

These files must be kept together in the same folder. The .ipynb file calls on the other files for functionality.

Using the Notebook with VS Code –

This tutorial assumes the use of sequences coming from the NADEL Library. If you need to use another library, see the **Customization.pdf**

Use VS Code to open **DelScreenProcessing.ipynb**



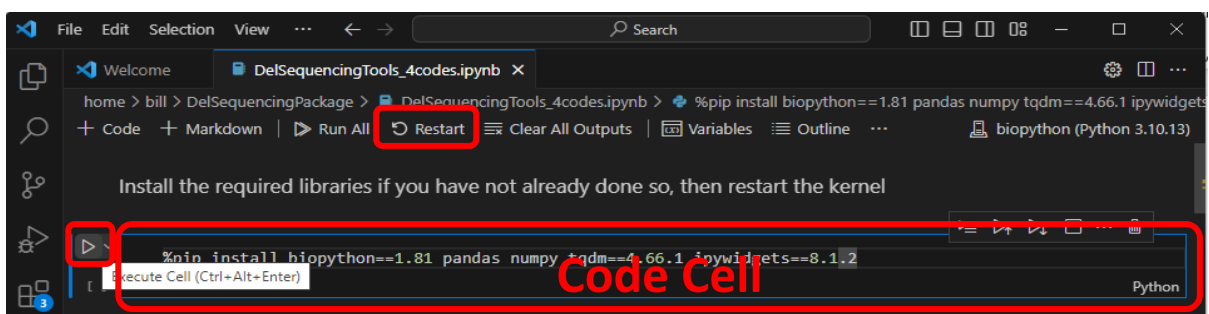
Now tell VS Code to use the biopython venv we created earlier to run this notebook (above screenshot).

In the upper right, click: Select Kernel > Python Environments... > biopython (Python 3.10.8)

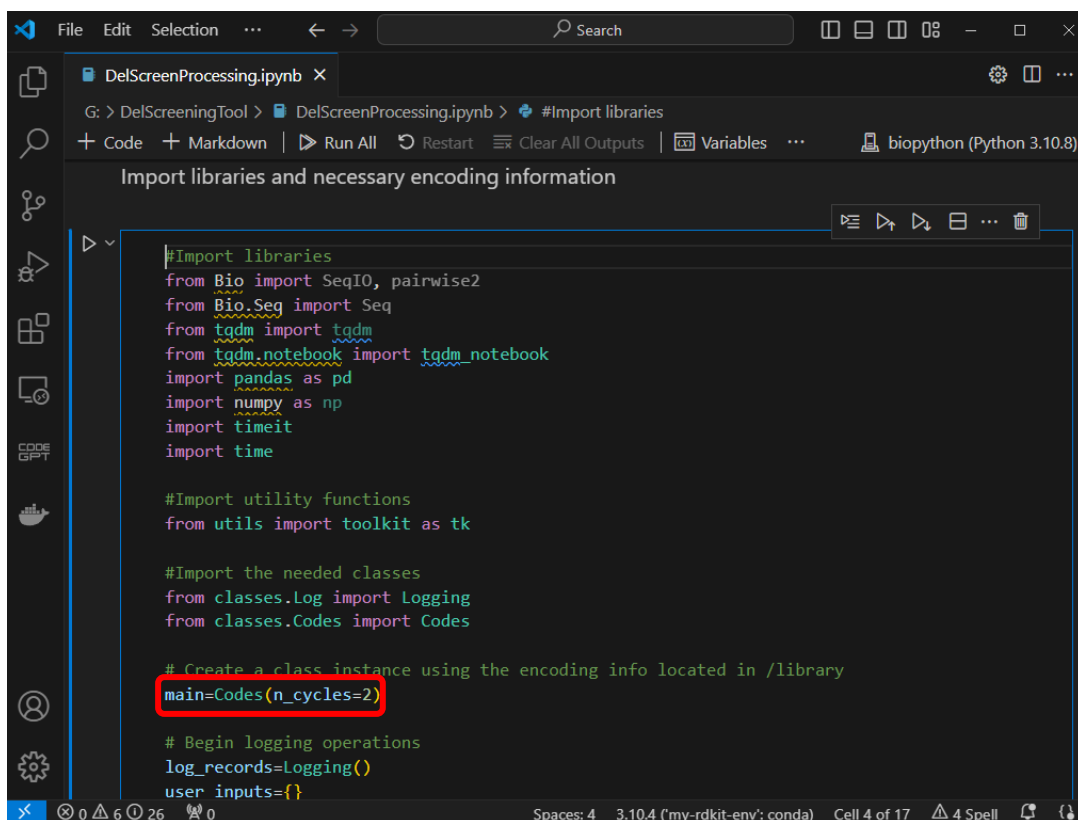
Notebook cells are run by mousing over them and then clicking the “Play” icon that appears just to the left of the cell.

For the first time running this notebook, lets go ahead and run the top cell to make sure all the necessary libraries are installed and up to date. After doing this restart the kernel. See the screenshot below for reference.

You should only need to run this topmost cell once.



Next, we want to import the python libraries we are going to be actively calling upon with the rest of code cells (below). We also will import the data contained in the other .py files and define a class (“main”, in this example) to hold our encoding sequence information.

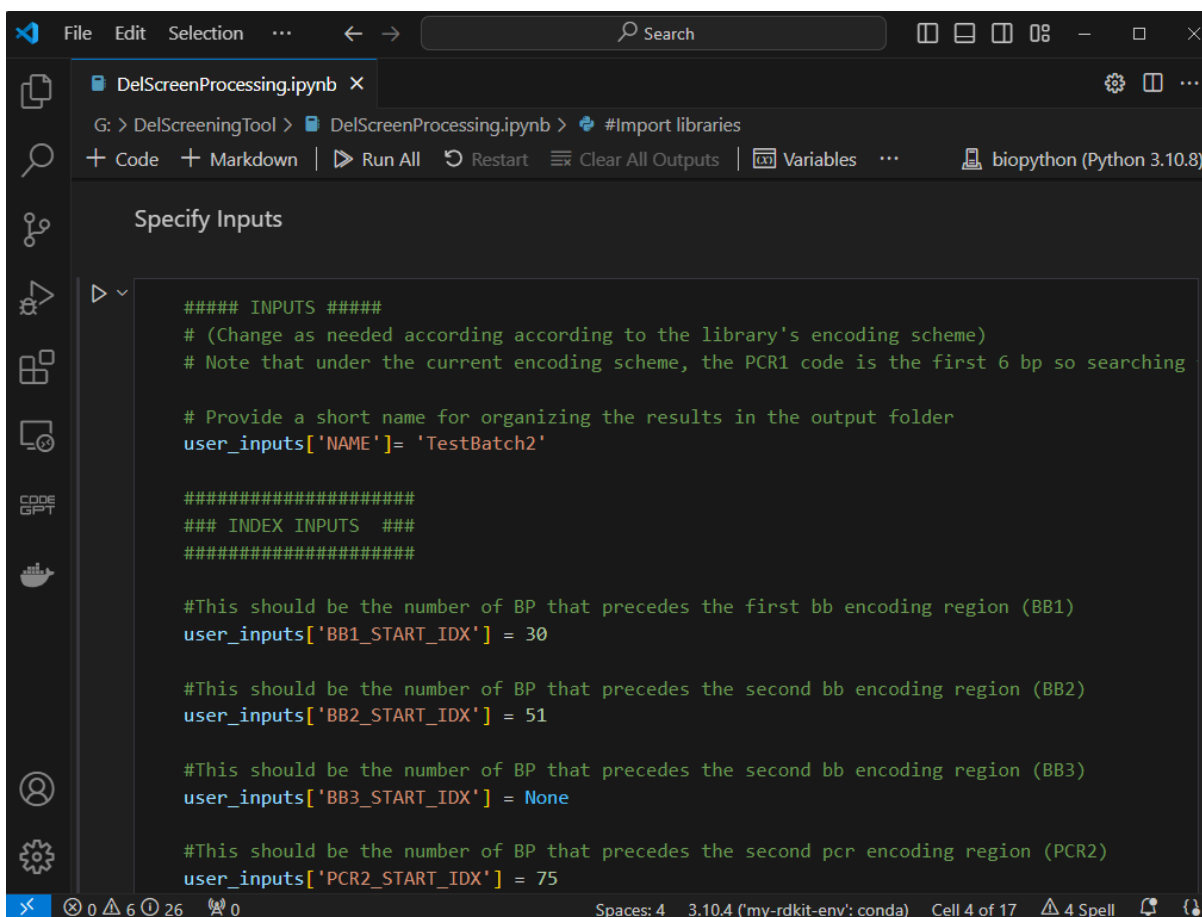


Be sure to set the “n_cycles” argument to the appropriate number of diversity elements in your library, then run the cell.

Before running the next cell which accepts all the rest of the user inputs (below), familiarize yourself with the different variables in this section. These are already predefined for the NADEL library but you will need to change them accordingly if a different library is being used. For more information see the **Customization.pdf**.

For now, if you are using a fastq file other than the current one provided for testing (ExampleData.fastq), drag the fastq file you want to process into the same folder with the rest of the package files and update the RAW_FASTQ_FILE variable to match the name of your file. Use the NAME variable to enter a short identifier for the screen(s) associated with your raw fastq data (this will be the name of the subfolder used for outputting results. For libraries other than NADEL, you will have to determine the BB and PCR index values.

After making any other desired changes, run the cell.



The screenshot shows a Jupyter Notebook window titled "DelScreenProcessing.ipynb". The interface includes a top menu bar with "File", "Edit", "Selection", and a search bar. Below the menu is a toolbar with icons for file operations and a "Run All" button. The main area displays a code cell with the following text:

```
##### INPUTS #####
# (Change as needed according to the library's encoding scheme)
# Note that under the current encoding scheme, the PCR1 code is the first 6 bp so searching

# Provide a short name for organizing the results in the output folder
user_inputs['NAME'] = 'TestBatch2'

#####
### INDEX INPUTS ###
#####

#This should be the number of BP that precedes the first bb encoding region (BB1)
user_inputs['BB1_START_IDX'] = 30

#This should be the number of BP that precedes the second bb encoding region (BB2)
user_inputs['BB2_START_IDX'] = 51

#This should be the number of BP that precedes the second bb encoding region (BB3)
user_inputs['BB3_START_IDX'] = None

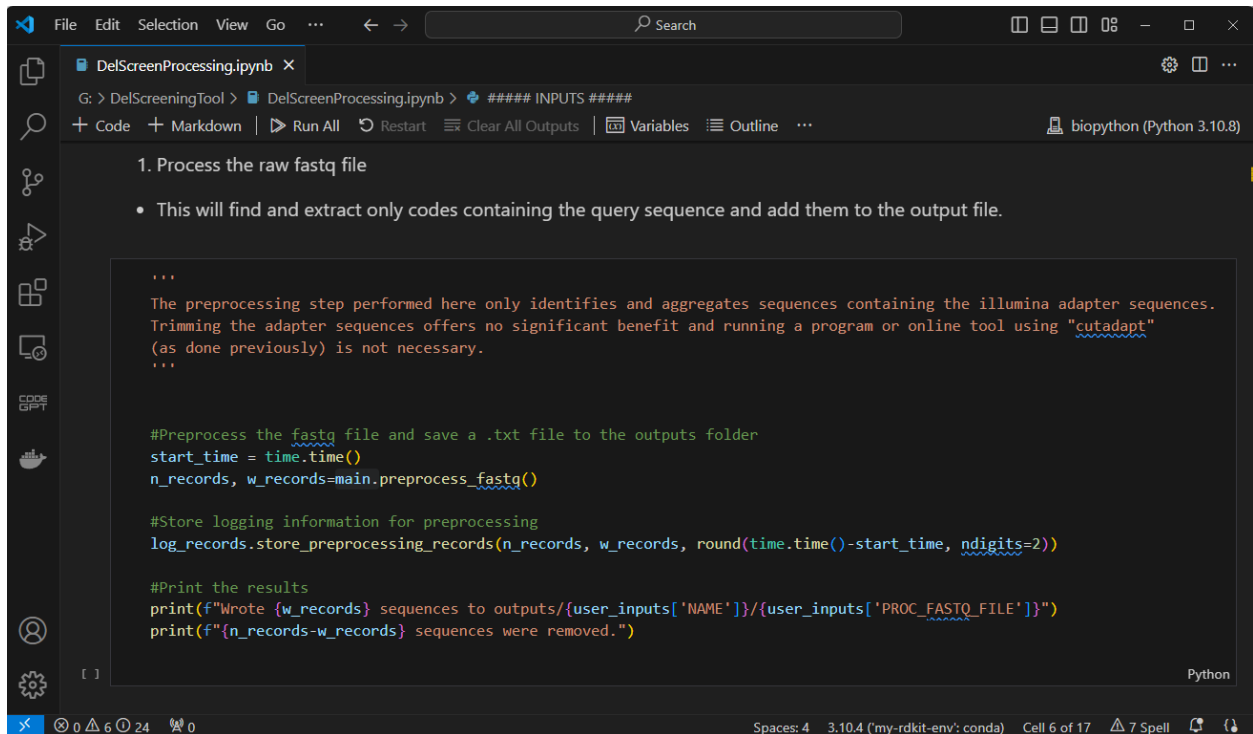
#This should be the number of BP that precedes the second pcr encoding region (PCR2)
user_inputs['PCR2_START_IDX'] = 75
```

The bottom status bar indicates "Spaces: 4", "3.10.4 ('my-rdkit-env': conda)", "Cell 4 of 17", and "4 Spell".

We are now ready to begin our Data processing process.

First let's run the next cell (below) to extract only sequence reads that contain an Illumina primer sequence and add the sequence information to a separate txt file (created automatically in the outputs

folder). For our purposes, the fastq IDs and qualities are removed since they only increase the file size and runtime.



The screenshot shows a Jupyter Notebook titled 'DelScreenProcessing.ipynb' in a web browser. The interface includes a top menu bar (File, Edit, Selection, View, Go, Search), a toolbar with icons for file operations, and a status bar at the bottom showing 'Spaces: 4', '3.10.4 (my-rdkit-env: conda)', 'Cell 6 of 17', and '7 Spell'. The notebook content is as follows:

```
##### INPUTS #####

1. Process the raw fastq file

• This will find and extract only codes containing the query sequence and add them to the output file.

...
The preprocessing step performed here only identifies and aggregates sequences containing the illumina adapter sequences.
Trimming the adapter sequences offers no significant benefit and running a program or online tool using "cutadapt"
(as done previously) is not necessary.
...

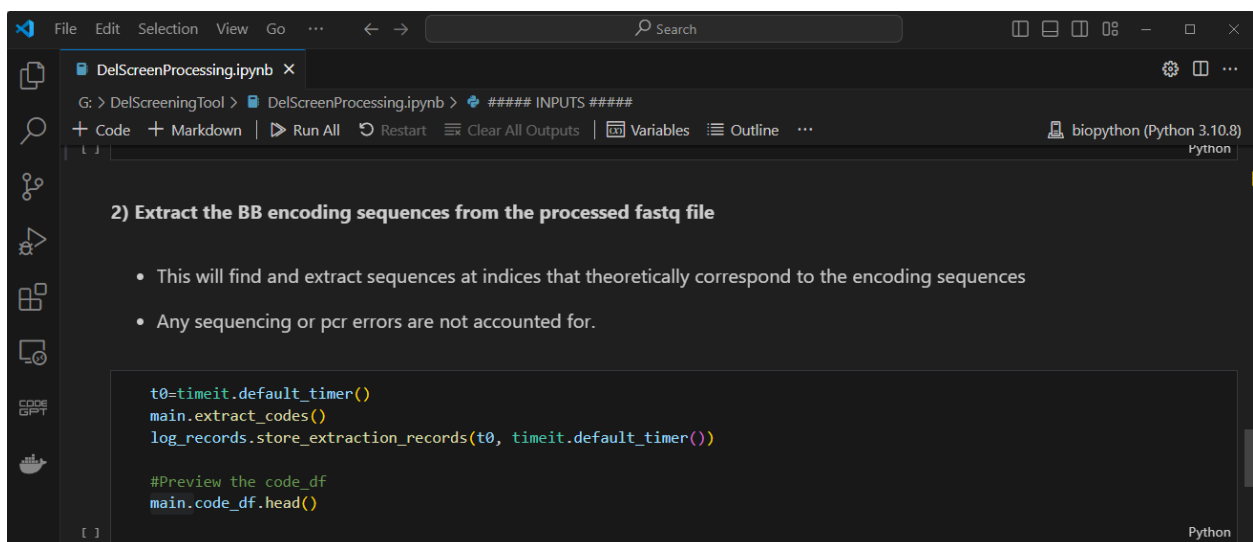
#Preprocess the fastq file and save a .txt file to the outputs folder
start_time = time.time()
n_records, w_records=main.preprocess_fastq()

#Store logging information for preprocessing
log_records.store_preprocessing_records(n_records, w_records, round(time.time()-start_time, ndigits=2))

#Print the results
print(f"Wrote {w_records} sequences to outputs/{user_inputs['NAME']}/{user_inputs['PROC_FASTQ_FILE']}")
print(f"{n_records-w_records} sequences were removed.")
```

Next, let's extract the encoding information from the preprocessed sequence reads.

This method "slices" the encoding sequences out of the sequence read using their expected positions.



The screenshot shows the same Jupyter Notebook interface, now displaying the second step of the pipeline. The content is as follows:

```
2) Extract the BB encoding sequences from the processed fastq file

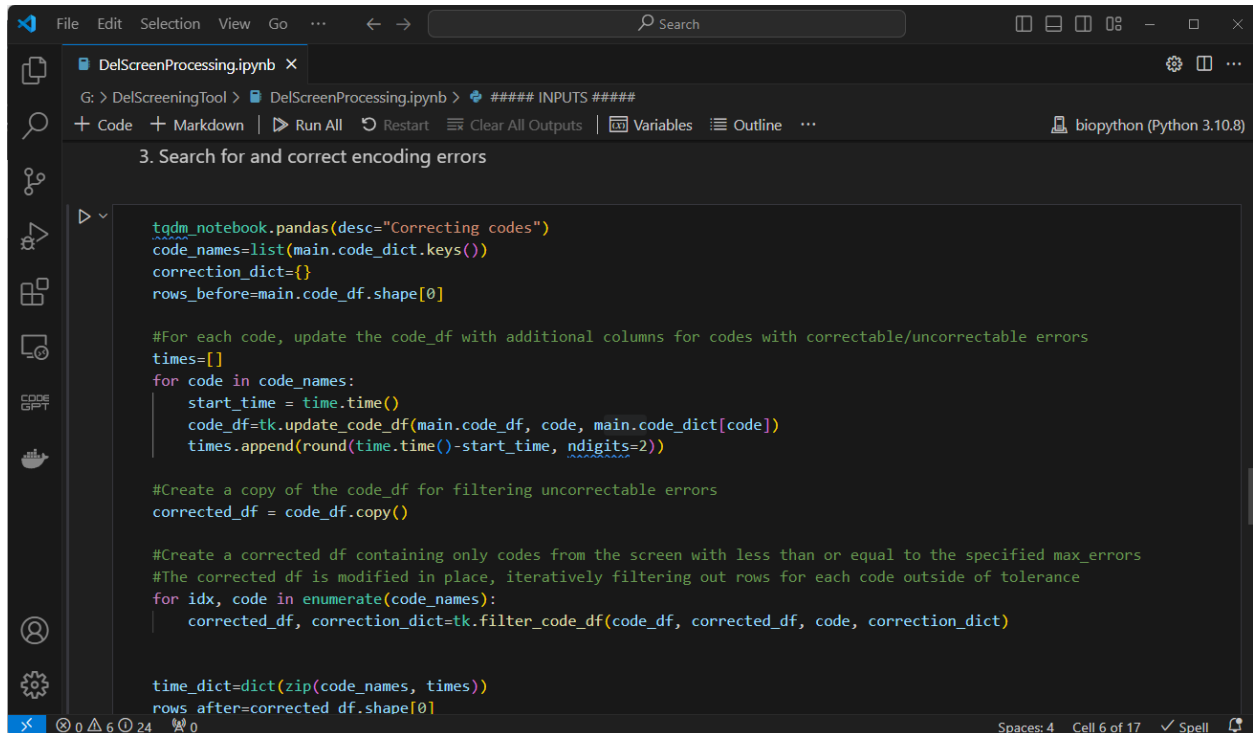
• This will find and extract sequences at indices that theoretically correspond to the encoding sequences
• Any sequencing or pcr errors are not accounted for.

t0=timeit.default_timer()
main.extract_codes()
log_records.store_extraction_records(t0, timeit.default_timer())

#Preview the code_df
main.code_df.head()
```

We next will apply a code correction step. Since our encoding sequences are unique from each other at 2 or more positions, if only a single error is present in the sequence, we can correct it in some cases. Reads with 2 or more errors are rejected and removed from the final counts.

Run the following cell (below) to perform the encoding corrections.



The screenshot shows a Jupyter Notebook interface with a dark theme. The active cell is titled "3. Search for and correct encoding errors". The code within the cell performs the following steps: it imports 'tqdm_notebook.pandas' with the description "Correcting codes"; it lists the keys of 'main.code_dict' as 'code_names'; it initializes an empty 'correction_dict' and gets the initial number of rows from 'main.code_df'; it enters a loop over 'code_names' where it updates 'code_df' for each code, recording the time taken; it creates a copy of 'code_df' named 'corrected_df'; it then iteratively filters 'corrected_df' for each code based on the 'correction_dict' using a function 'tk.filter_code_df'; finally, it creates a 'time_dict' from the recorded times and updates the number of rows in 'corrected_df'.

```
tqdm_notebook.pandas(desc="Correcting codes")
code_names=list(main.code_dict.keys())
correction_dict={}
rows_before=main.code_df.shape[0]

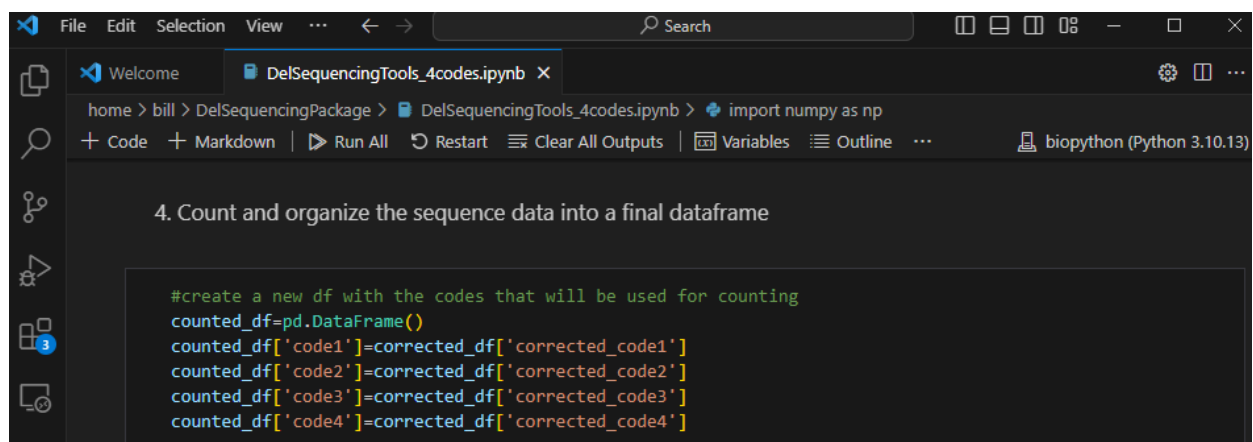
#For each code, update the code_df with additional columns for codes with correctable/uncorrectable errors
times=[]
for code in code_names:
    start_time = time.time()
    code_df=tk.update_code_df(main.code_df, code, main.code_dict[code])
    times.append(round(time.time()-start_time, ndigits=2))

#Create a copy of the code_df for filtering uncorrectable errors
corrected_df = code_df.copy()

#Create a corrected df containing only codes from the screen with less than or equal to the specified max_errors
#The corrected df is modified in place, iteratively filtering out rows for each code outside of tolerance
for idx, code in enumerate(code_names):
    corrected_df, correction_dict=tk.filter_code_df(code_df, corrected_df, code, correction_dict)

time_dict=dict(zip(code_names, times))
rows_after=corrected_df.shape[0]
```

Next, we will combine and count identical reads into a final table which is done by running the penultimate cell (below).



The screenshot shows a Jupyter Notebook interface with a dark theme. The active cell is titled "4. Count and organize the sequence data into a final dataframe". The code within the cell creates a new 'counted_df' as a 'pd.DataFrame()' and then maps columns from 'corrected_df' to 'counted_df' for 'code1' through 'code4'.

```
#create a new df with the codes that will be used for counting
counted_df=pd.DataFrame()
counted_df['code1']=corrected_df['corrected_code1']
counted_df['code2']=corrected_df['corrected_code2']
counted_df['code3']=corrected_df['corrected_code3']
counted_df['code4']=corrected_df['corrected_code4']
```

Finally, run the last cell (below) to output your results as two csv files. One csv will have absolute sequence counts and the other will have normalized sequence counts.

DelScreenProcessing.ipynb

G: > DelScreeningTool > DelScreenProcessing.ipynb > ##### INPUTS #####

+ Code + Markdown ▶ Run All ↺ Restart 🗑 Clear All Outputs 📄 Variables 📖 Outline ...

biopython (Python 3.10.8)

5. Output the results in two csv files:

- 1. As absolute sequence counts
- 2. As normalized sequence counts

```
# Merge column names
merged_columns = [''.join(col).strip() for col in final_df.columns.values]
merged_columns = [str(col).replace('1a-', '') for col in merged_columns]
merged_columns = [str(col).replace('1b-', '') for col in merged_columns]

# Assign merged column names to DataFrame
final_df.columns = merged_columns

final_df.to_csv(f'{main.output_location}/Results_AbsSeqCounts.csv')

#Convert to normalized sequence counts
#pcr_list=list(pcr_primer_1.values())
for pcr in final_df.columns[4:]:
    sum_value = final_df[pcr].sum()
    final_df[pcr]=final_df[pcr].apply(lambda count: (count/sum_value)*(len(final_df)))

final_df.to_csv(f'{main.output_location}/Results_NormSeqCounts.csv')
```

0 0 6 0 24 0

Spaces: 4 Cell 6 of 17 ✓ Spell