

CAVA All Hands Meeting

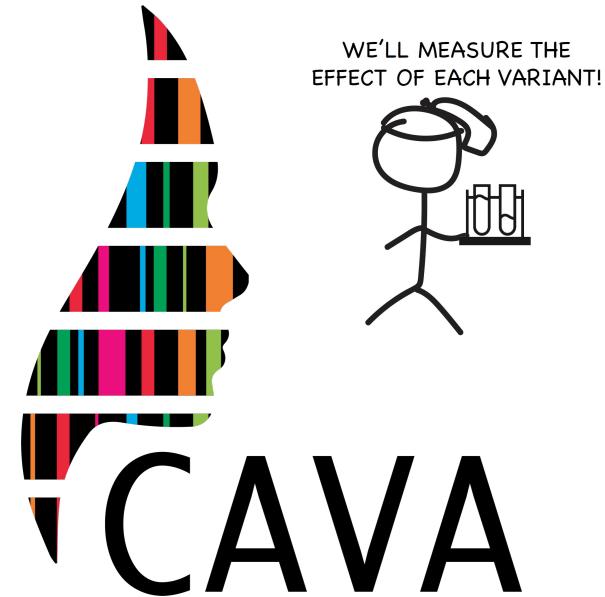
Tues, April 30th 2024

4pm-5pm Pacific | 11am-12pm Melbourne

Location: Foege S448 and zoom

## AGENDA:

CountESS with Nick Moore!



We build high resolution  
sequence-function maps for  
clinically actionable genes



# CountESS

Count-based Experimental Scoring and Statistics



Hollerith 1890 Tabulating Machine and Sorting Box

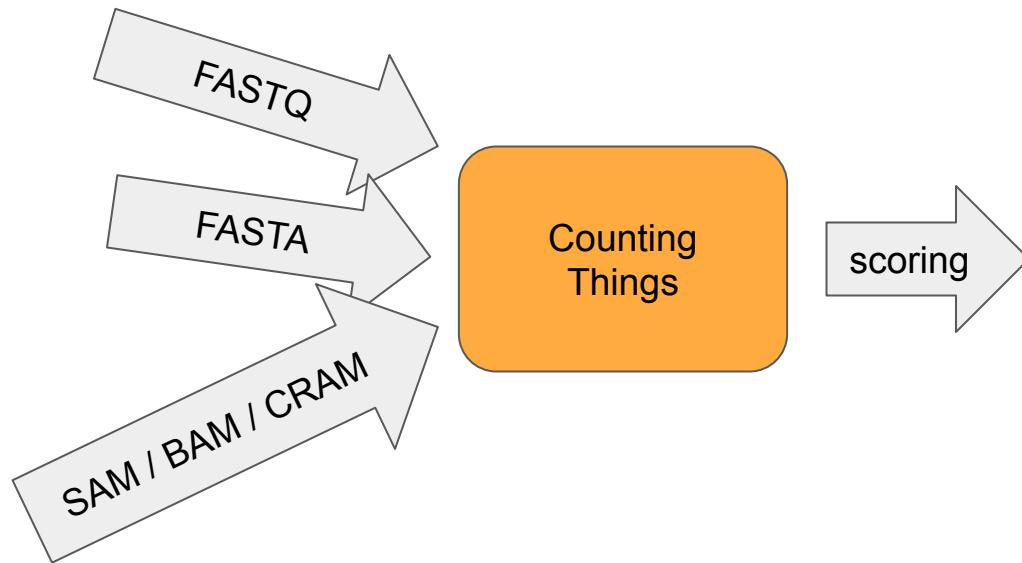
*Image: Adam Schuster, CC BY 2.0, via Wikimedia Commons*

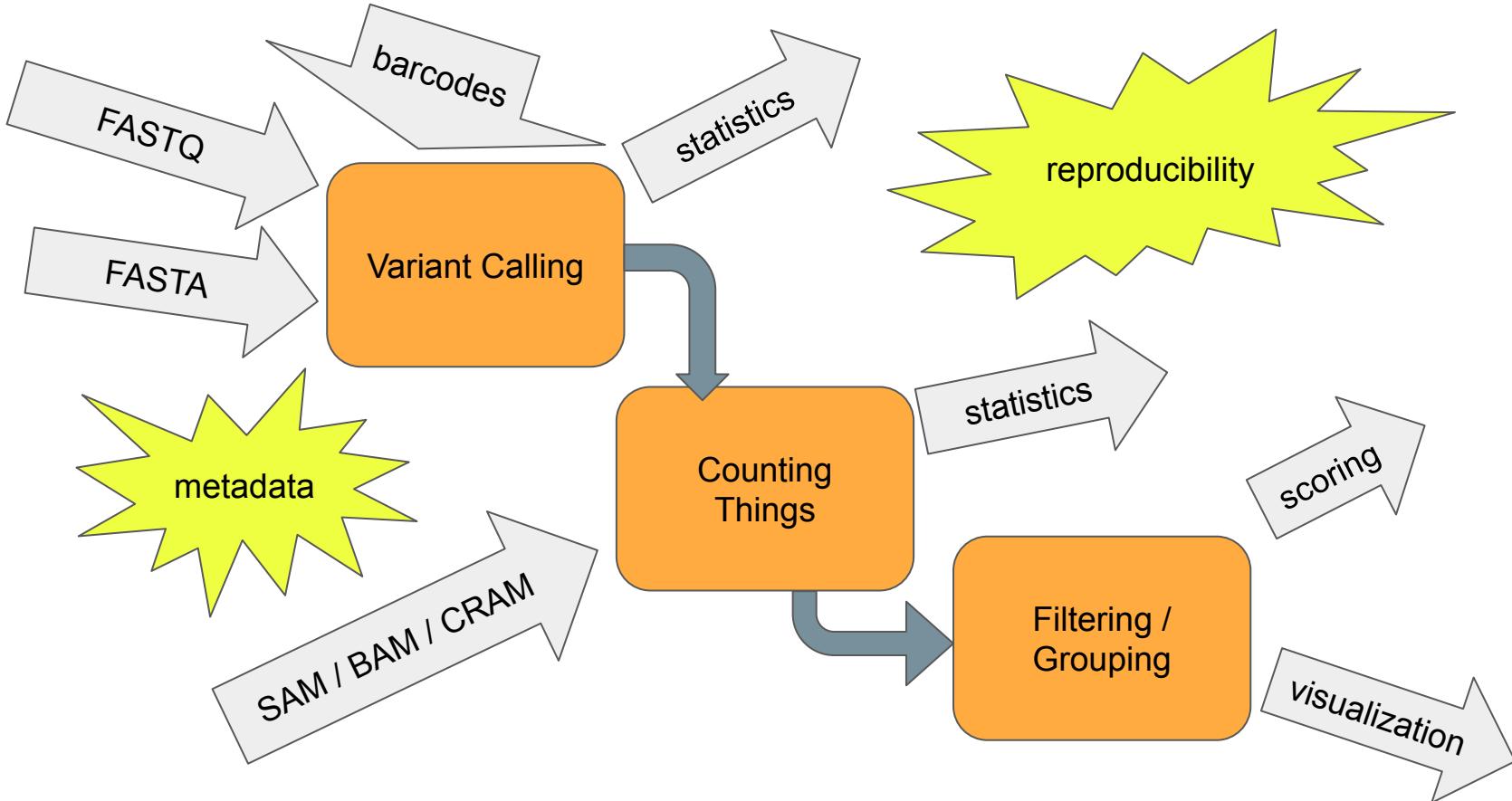
<https://commons.wikimedia.org/wiki/File:HollerithMachine.CHM.jpg>

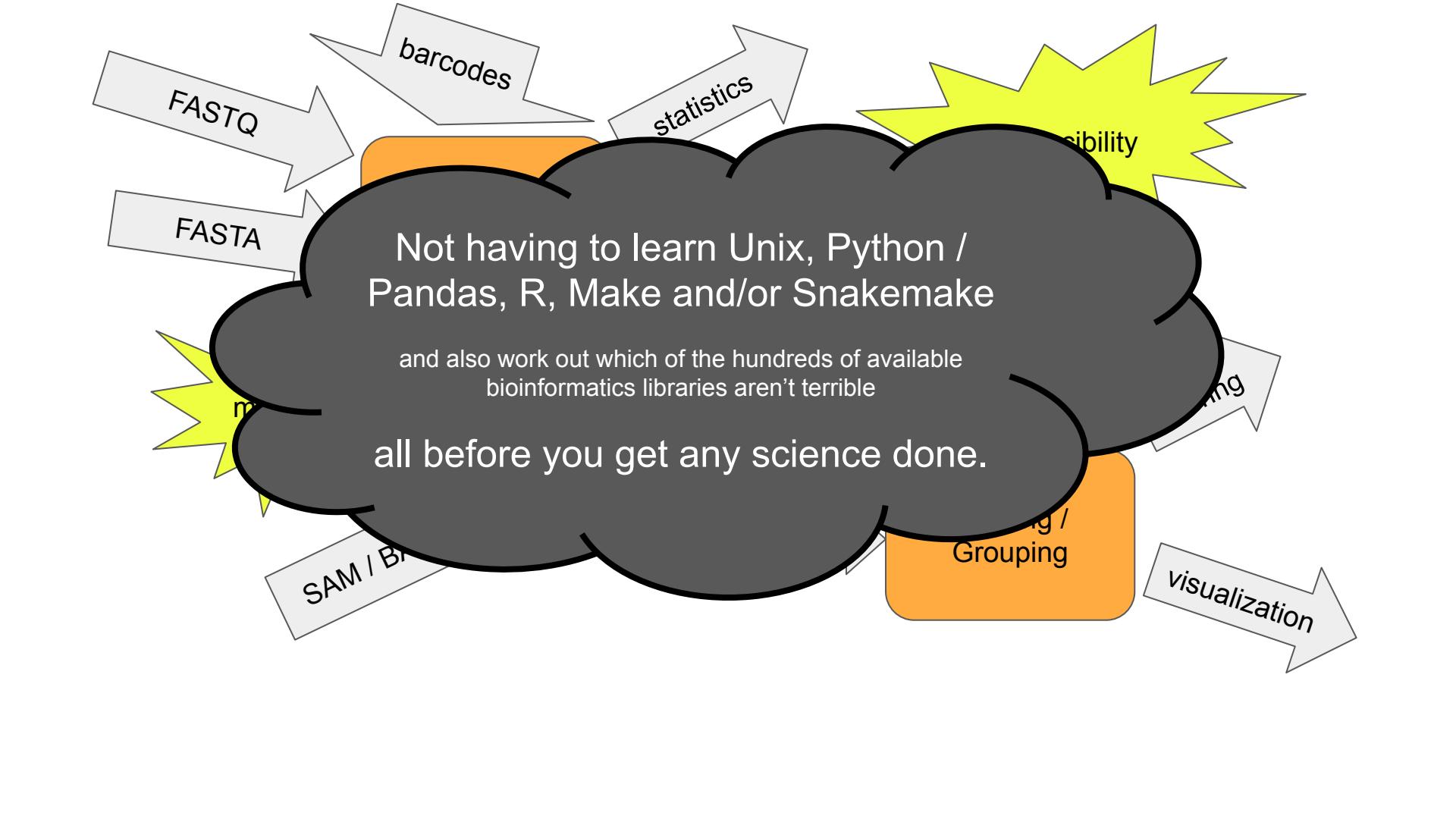
Counting  
Things

```
sort sequences.txt  
| uniq -cd
```

Counting  
Things







Not having to learn Unix, Python /  
Pandas, R, Make and/or Snakemake

and also work out which of the hundreds of available  
bioinformatics libraries aren't terrible

all before you get any science done.

FASTQ

barcodes

statistics

FASTA

sibility

m

SAM / BAM

ing / Grouping

Visualization



CrossMark

# Gene name errors are widespread in the scientific literature

Mark Ziemann<sup>1</sup>, Yotam Eren<sup>1,2</sup> and Assam El-Osta<sup>1,3\*</sup>

## Abstract

The spreadsheet software Microsoft Excel, when used with default settings, is known to convert gene names to dates and floating-point numbers. A programmatic scan of leading genomics journals reveals that approximately one-fifth of papers with supplementary Excel gene lists contain erroneous gene name conversions.

**Keywords:** Microsoft Excel, Gene symbol, Supplementary data

**Abbreviations:** GEO, Gene Expression Omnibus; JIF, journal impact factor

The problem of Excel software (Microsoft Corp., Redmond, WA, USA) inadvertently converting gene symbols to dates and floating-point numbers was originally described in 2004 [1]. For example, gene symbols such as *SEPT2* (Septin 2) and *MARCH1* [Membrane-Associated Ring Finger (C3HC4) 1, E3 Ubiquitin Protein Ligase] are converted by default to '2-Sep' and '1-Mar', respectively. Furthermore, RIKEN identifiers were described to be automatically converted to floating point numbers (i.e. from accession '2310009E13' to '2.31E+13'). Since that report, we have uncovered further instances where gene symbols were converted to dates in supplementary data of recently published papers (e.g. '*SEPT2*' converted to '2006/09/02').

## Enrich2:

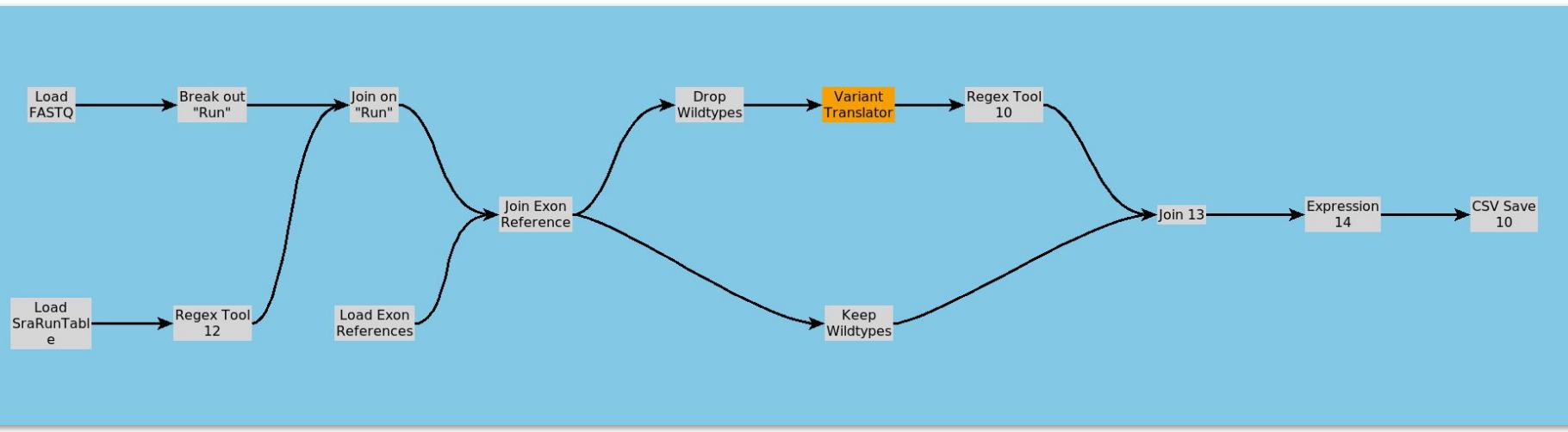
- GUI based
- Python 2 - EOL!
- Specific to a few experiments
- Difficult to extend

## CountESS:

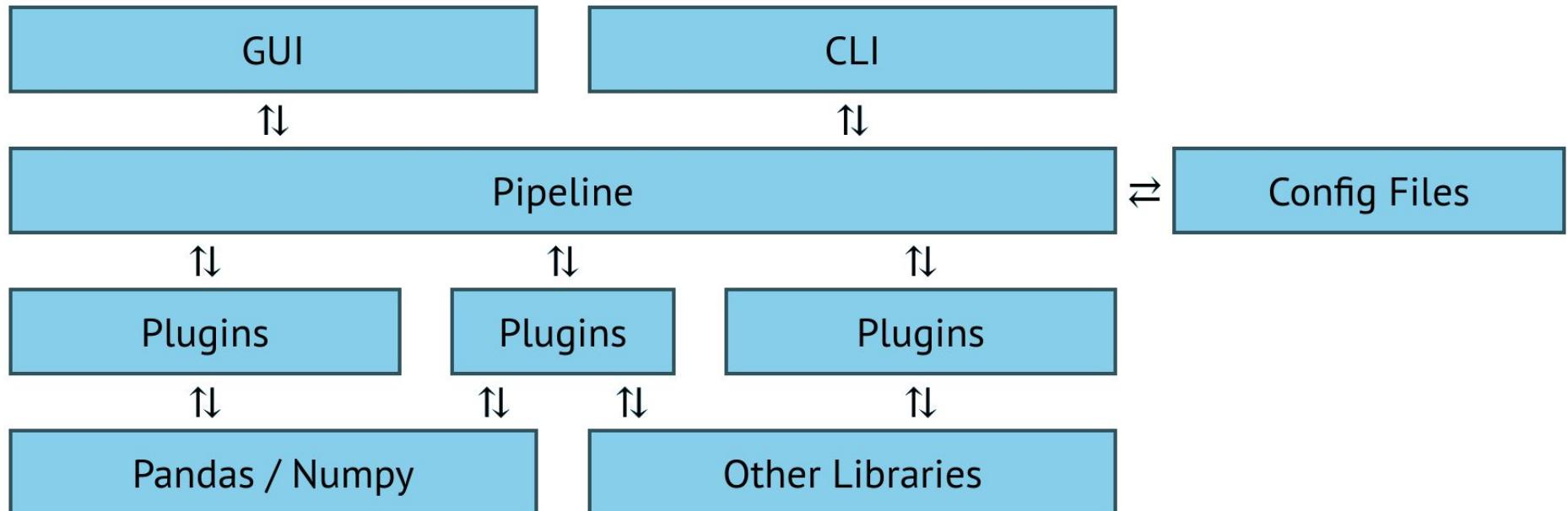
- GUI based
- Python 3.9+
- Plugin architecture
- Simple, single text config file
- Flexible pipeline

# CountESS Plugins & Pipelines

- Some plugins included, plus it's easy to install / write new ones.
- Plugins are connected together in a potentially complex pipeline.



# CountESS Architecture



# CountESS Documentation

<https://countess-project.github.io/CountESS/>

# Installing CountESS

```
pip install countess
```

# Running CountESS

countess\_gui

countess\_cmd

## NEW

### Choose Plugin

CSV Load

Loads data from CSV or similar delimited text files and assigns types to columns

DataTable

enter small amounts of data directly

FASTQ Load

Loads counts from FASTQ files containing either variant or barcodes

Mutagenize

Provides all mutations of a sequence

Regex Reader

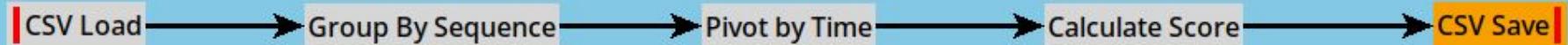
Loads arbitrary data from line-delimited files

NEW

# Examples

<https://github.com/CountESS-Project/countess-demo/>

# Example 1: Reading & Counting Sequences



sequence, time

AGTCGAGGACATGGTGAGT , 1  
GATGTCCTAACGGTCGATT , 1  
TTCAGTACCTAAACTATGTT , 1  
TTTATATTTCGAGTGAATGT , 1  
CAACGAGAGATGTAGGAGAA , 1  
GTAACAGGAGTCATGTTCC , 1  
(etc)

sequence, time

GTTCGGCTACCAGGCGAGGA , 2  
TCTCCTTGACGACTCGCAC , 2  
GTTAAGGGCCTCCAAGGAGA , 2  
AGAGACTGCGCACCGCCCCGC , 2  
GTACACTCAATAAGTACTGT , 2  
TCATTATACTTCTCAATACT , 2  
(etc)

**CSV Load**

CSV Load 0.0.55 — Loads data from CSV or similar delimited text files and assigns types to columns

[add notes](#)

**Files** [+](#)

	Filename	
File	sequences_1.csv	<input checked="" type="checkbox"/>
File	sequences_2.csv	<input checked="" type="checkbox"/>

Delimiter ,

Quoting None

Comment None

CSV file has header row?

Filename Column

**Columns** [+](#)

Column	Column Name	Column Type	Index?
sequence	string	<input type="checkbox"/>	<input checked="" type="checkbox"/>
time	integer	<input type="checkbox"/>	<input checked="" type="checkbox"/>

Dataframe Preview 86374 rows

	sequence	time
1	TGCTTGATGGATAGCTCAAAG	1
1	TCTTGGGTTTGATCTTTTC	1
1	AACATTGTATATTCAAGCGTA	1
1	CAACCTCAGGGAGCCTTAGC	1
1	CTATACTCAGAGTACGACT	1
1	CCTAGGTGGATAGAACGCAA	1
1	TCTCACTCGAGCTTGAGCA	1
1	CGCGCTTGCATATGCTCAT	1
1	TGAATAGAAAGCGCATACCG	1

```

graph TD
    A[CSV Load] --> B[Group By Sequence]
    B --> C[Pivot by Time]
    C --> D[Calculate Score]
    D --> E[CSV Save]
  
```

**Group By Sequence**

Group By 0.0.55 — Group records by column(s) and calculate aggregates

[add notes](#) [i](#)

**Columns**

	Index	Count	Min	Max	Sum	Mean
"sequence"	<input checked="" type="checkbox"/>	<input type="checkbox"/>				
"time"	<input checked="" type="checkbox"/>	<input type="checkbox"/>				

Join Back?

**Dataframe Preview 1999 rows**

sequence object ▼	time int64 ▲	count int64 ▲
TTTAAGCTGGACTAGATGC	2	46
TTTAAGCTGGACTAGATGC	1	67
TTTAGTTATCGTGTGTTA	2	46
TTTGTTAGTTATCGTGTGTTA	1	64
TTTGTAAAGAATGACCAACCG	2	12
TTTGTAAAGAATGACCAACCG	1	16
TTTGGCATACTGTAATCCCAT	2	81
TTTGGCATACTGTAATCCCAT	1	99
TTTGAGCTCGATTAGAA	2	15

```

graph TD
    A[CSV Load] --> B[Group By Sequence]
    B --> C[Pivot by Time]
    C --> D[Calculate Score]
    D --> E[CSV Save]
  
```

sequence	time	count
AAAAA	1	5
AAAAA	2	7
CCCCC	1	8
GGGGG	1	10
GGGGG	2	3



sequence	count_time_1	count_time_2
AAAAA	5	7
CCCCC	8	0
GGGGG	10	3

## Pivot by Time



Pivot Tool 0.0.55 — Groups a dataframe and pivots column values into columns.

Expanded column values are duplicated for each combination of pivot values. Missing values default to 0, and duplicate values are summed.

[add notes](#)

"sequence"

Index

"time"

Pivot

"count"

Expand

Dataframe Preview 1000 rows

sequence object (index)	count_time_1 int64	count_time_2 int64
AAAAATCCGTAGGGGTTGCC	35	25
AAAATTTGAAGTGGGTACG	19	16
AAAAGAAGCTCTAGTATATT	96	71
AAAATAGAACCGTGGCACCT	29	22
AAACACTGGTTAGACCCAAG	88	65
AAACAGGTGTCGCCCAAGC	25	18
AAACCTTGGGGACACCCGGC	49	36
AAACTTTCAGTCTCAGGAGA	7	5

## Calculate Score

CSV Load

Group By Sequence

Pivot by Time

Calculate Score

CSV Save

Python Code 0.0.55 — Apply python code to each row.

Columns are mapped to local variables and back. If you assign to a variable called "`_filter`", only rows where that value is true will be kept.

[add notes](#)

```
score = (
    count_time_2 / count_time_1
    if count_time_1 else None
)
```

Python Code

Drop Null Columns?



Dataframe Preview 1000 rows

sequence	count_time_1	count_time_2	score
object	int64	int64	float64
AAAAATCCGTAGGGGTTGCC	35	25	0.714285714286
AAAACTTGAAGTGGGTACG	19	16	0.842105263158
AAAAAGAAGCTCTAGTATATT	96	71	0.739583333333
AAAATAGAACCGTGGCACCT	29	22	0.758620689655
AAACACTGGTTAGACCCAAG	88	65	0.738636363636
AAACAGGTGTCGCCCAAGC	25	18	0.72
AAACCTTGGGGACACCCGGC	49	36	0.734693877551
AAACTTTCAGTCTCAGGAGA	7	5	0.714285714286

## CSV Save

CSV Save 0.0.55 — Save data as CSV or similar delimited text files

[add notes](#)

CSV header row?

Filename Delimiter Quote all Strings **CSV Load****Group By Sequence****Pivot by Time****Calculate Score****CSV Save**

Text Preview 1001 Lines

```
sequence,count_time_1,count_time_2,score
AAAAATCCGTAGGGTTGCC,35,25,0.7142857142857143
AAAACTTTGAAGTGGGTACG,19,16,0.8421052631578947
AAAAGAAGCTCTAGTATATT,96,71,0.7395833333333334
AAAATAGAACCGTGGCACCT,29,22,0.7586206896551724
AAACACTGGTTAGACCAAG,88,65,0.7386363636363636
AACACGGTGTGCCAAGC,25,18,0.72
AACACCTTGGGGACACCCGGC,49,36,0.7346938775510204
AAACTTCAGTCTCAGGAGA,7,5,0.7142857142857143
AAAGATTAACACATAAGCTAA,90,66,0.7333333333333333
AAAGCTAACATAAGCTAA,17,14,0.8235294117647058
AAAGCTAACATAACGGCAGATCT,12,9,0.75
AAAGCTAGATACGAGGACCT,69,57,0.8260869565217391
```

## Example 2: Translating Barcodes & Calling Variants

barcode, sequence

ATTCCCGTAATCTACGATTA, ATGCTTGTACGGTGGTGCCTGGCTTATCTATCTAGAT  
CCGTCTCCGAGTCACGGTCGAATTAGGTACTGCACTATCCTTGAGGCCGGAAAGGGCAC  
AAGGGCCGACCCTTGTGGATAAAATTGCTAAGAGGAAGGTCTAG  
**TTACGGTCTGCGTTGGAATC**, ATGCTTGTACGGTGGTGCCTGGCTTATCTATCTAGAT  
CCGTCTCCGAGTCACGGTCGAATTAGGTACTGCACTATCCTTGAGGCCGGAAAGGGCAC  
AAGGGCCGACCCTTGTGGATAAAATTGCTAAGAGGAAGGTCTAG  
**AGGGCCGTGCCAAGTGCAGT**, ATGCTTGTACGGTGGTGCCTGGCTTATCTATCTAGAT  
CCGTCTCCGAGTCACGGTCGAATTAGGTACTGCACTATCCTTGAGGCCGGAAAGG**A**CCAC  
AAGGGCCGACCCTTGTGGATAAAATTGCTAAGAGGAAGGTCTAG  
**TGTAGTGCCGTATTGTGGC**, ATGCTTGTACGGTGGTGCCTGGCTTATCTATCTAGAT  
CCGTCTCCGAGTCACGGTCGAATTAGGTACTGCACTATCCTTGAGGC**A**GGAAAGGGCAC  
AAGGGCCGACCCTTGTGGATAAAATTGCTAAGAGGAAGGTCTAG

**CSV Load Sequences**

CSV Load 0.0.55 — Loads data from CSV or similar delimited text files and assigns types to columns

Files + add notes

File	Filename	sequences_1.csv	X
File	Filename	sequences_2.csv	X
Delimiter	,	New Config	Load Config
Quoting	None	Save Config	Export Config
Comment	None	Tidy Graph	Run

CSV file has header row?

Filename Column

Columns +

Column Name	barcode	Index
Column	barcode	✓
Column	time	✓

barcode object

```

TGCTTGTGGATAGCTCAAG
TCTTGGGTTGCATCTTTC
AACATTGTATATTCAAGCTA
CAACCTCCGGAGCCTTAGC
CTATACTCAGAGTAAGACT
CCTAGGTGGATAGAACGAA
TCTCACTGGAGCTTGAAGCA
CGCGCTTGCATATGCTCAT
TGAATAGAAAGGCCATACCG

```

**Group By Barcode**

Group By 0.0.55 — Group records by column(s) and calculate aggregates

Columns

Index	Count	Min	Max	Sum	Mean
"barcode"	✓	X	X	X	X
"time"	✓	X	X	X	X

Join Back? X

Dataframe Preview 1999 rows

barcode object (index)	time int64 (index)	count int64
AAAAATCCGTAGGGTTGCC	1	35
AAAATCCGTAGGGTTGCC	2	25
AAAACTTGAAGTGGGTACG	1	19
AAAACTTGAAGTGGGTACG	2	16
AAAAGAAGCTAGTATATT	1	96
AAAAGAAGCTAGTATATT	2	71
AAAATAGAACCGTGGCACCT	1	29
AAAATAGAACCGTGGCACCT	2	22
AAACACTGGTTAGACCAAG	1	88

### CSV Load Barcodes

CSV Load 0.0.55 — Loads data from CSV or similar delimited text files and assigns types to columns

[add notes](#)

**Files** [+](#)

File	Filename
	barcodes.csv <a href="#">X</a>

Delimiter: ,  
Quoting: None  
Comment: None

CSV file has header row?

**Filename Column**

**Columns** [+](#)

Column	Column Name	Column Type	Index?
barcode	barcode	string	<input type="checkbox"/>
sequence	sequence	string	<input type="checkbox"/>

**Dataframe Preview 1000 rows**

barcode	sequence
object	object
AAAAATCCGTAGGGGTTGCC	ATGCTTGACGGGTGGTGCCTTATCTATCTAGATCCGTCCTCGAGTCACGGTC
AAAATTTGAAGTGGTAGC	ATGCTTGACGGGTGGTGCCTTATCTATCTAGATCCGTCCTCGAGTCACGGTC
AAAAGAAGCTCTAGTATATT	ATGCTTGACGGGTGGTGCCTTATCTATCTAGATCCGTCCTCGAGTCACGGTC
AAAATAGAACCGTGGCACCT	ATGCTTGACGGGTGGTGCCTTATCTATCTAGATCCGTCCTCGAGTCACGGTC
AAACACTGGTTAGACCCAAG	ATGCTTGACGGGTGGTGCCTTATCTATCTAGATCCGTCCTCGAGTCACGGTC
AAACAGGTGTCGCCCAAGC	ATGCTTGACGGGTGGTGCCTTATCTATCTAGATCCGTCCTCGAGTCACGGTC
AAACCTTGGGGACACCCGGC	ATGCTTGACGGGTGGTGCCTTATCTATCTAGATCCGTCCTCGAGTCACGGTC
AAACTTCAAGTCAGGAGA	ATGCTTGACGGGTGGTGCCTTATCTATCTAGATCCGTCCTCGAGTCACGGTC
AAAGATTAAAAATCTCAAT	ATGCTTGACGGGTGGTGCCTTATCTATCTAGATCCGTCCTCGAGTCACGGTC

```

graph TD
    A[CSV Load Sequences] --> B[Group By Barcode]
    C[CSV Load Barcodes] --> B
    B --> D[Join on Barcode]
    D --> E[Variant Translator]
    E --> F[Pivot by Time DNA]
    E --> G[Pivot by Time Protein]
    F --> H[Calculate Score DNA]
    F --> I[Calculate Score Protein]
    H --> J[CSV Save DNA]
    I --> K[CSV Save Protein]
  
```

**Join on Barcode**

Join 0.0.55 — Joins two Pandas Dataframes by indexes or columns i

**Input**

Join On — INDEX —

Required

Drop Column

**Input**

Join On — INDEX —

Required

Drop Column

**Dataframe Preview 1999 rows**

barcode object	time int64	count int64	sequence object
ACAGGTATGAGCTTGTA	1	53	TTGCTTGTACGGGTGGTGCCTGGCTTA
GTGAGTGTATTTGAGACT	1	96	TTGCTTGTACGGGTGGTGCCTGGCTTA
GTGAGTGTATTTGAGACT	2	74	TTGCTTGTACGGGTGGTGCCTGGCTTA
ACAGGTATGAGCTTGTA	2	39	TTGCTTGTACGGGTGGTGCCTGGCTTA
TAGGTTCAAATGAAAGATG	1	58	GTGCTTGTACGGGTGGTGCCTGGCTTA
TCATTATACTTCTCAACT	1	49	GTGCTTGTACGGGTGGTGCCTGGCTTA
TAGGTTCAAATGAAAGATG	2	38	GTGCTTGTACGGGTGGTGCCTGGCTTA
TCATTATACTTCTCAACT	2	40	GTGCTTGTACGGGTGGTGCCTGGCTTA
CGATGAGAAGGTAGTAAGAT	2	31	CTGCTTGTACGGGTGGTGCCTGGCTTA

```

graph TD
    A[CSV Load Sequences] --> B[Group By Barcode]
    C[CSV Load Barcodes] --> B
    B --> D[Join on Barcode]
    D --> E[Variant Translator]
    E --> F[Pivot by Time DNA]
    E --> G[Pivot by Time Protein]
    F --> H[Calculate Score DNA]
    F --> I[Calculate Score Protein]
    H --> J[CSV Save DNA]
    I --> K[CSV Save Protein]
  
```

## Variant Translator

Variant Translator 0.0.55 — Turns a DNA sequence into a HGVS variant code


[add notes](#)

Input Column

Reference Column

\*OR\* Reference Sequence

Output Column

Max Mutations

Protein Column

Protein Offset

Max Protein Variations

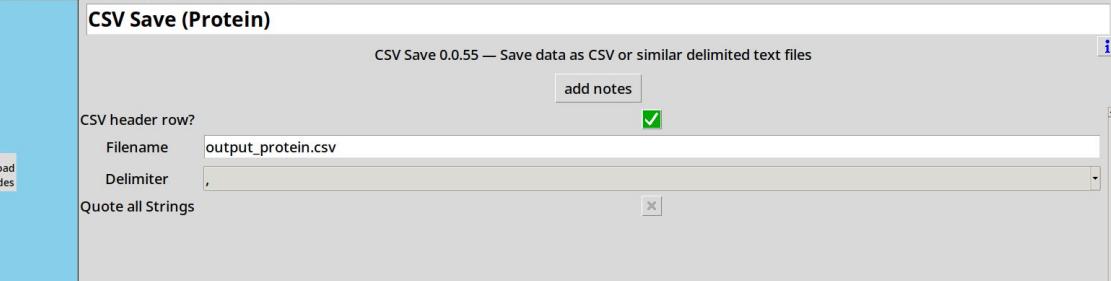
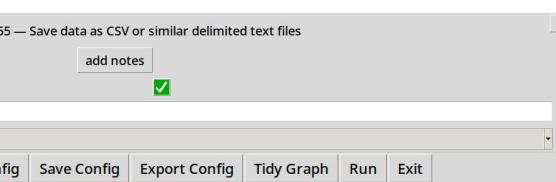
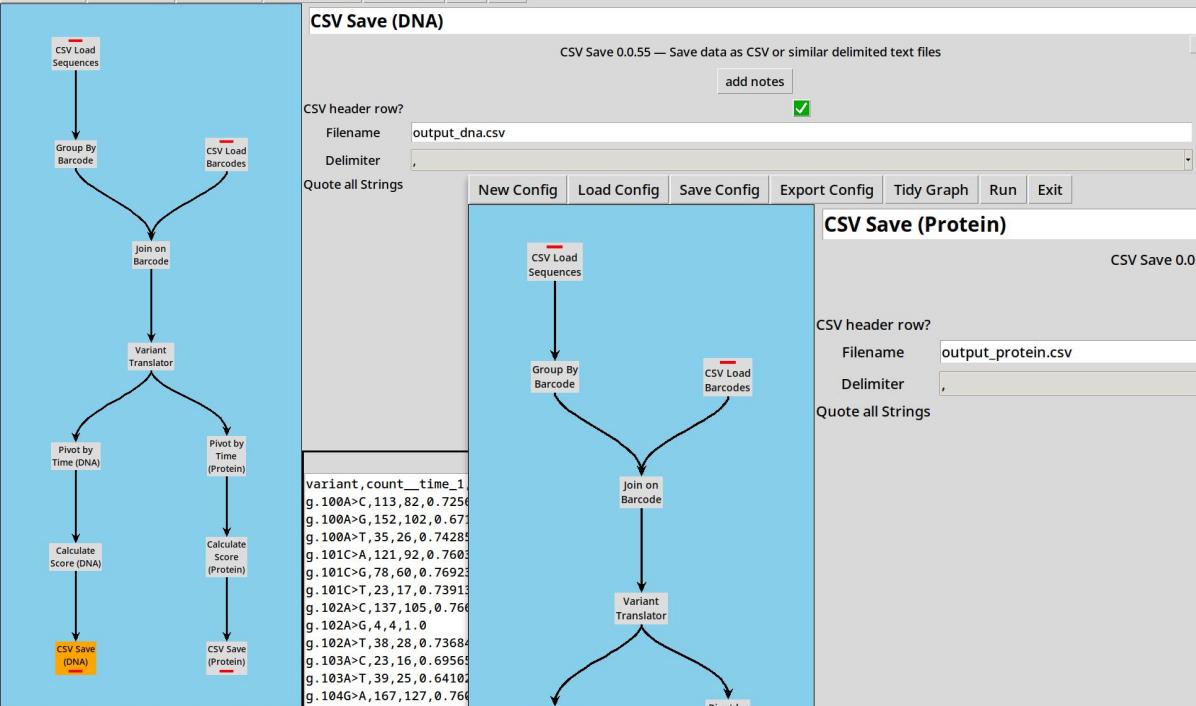
Drop unidentified variants

Drop Input Column(s)

Pivot by Time (DNA)  
Calculate Score (DNA)  
CSV Save (DNA)

Pivot by Time (Protein)  
Calculate Score (Protein)  
CSV Save (Protein)

Dataframe Preview 1999 rows					
barcode object	time int64	count int64	sequence object	variant object	protein object
AAAATCCGTAGGGTTGC	1	35	ATGCTTTGTACGGGTGGT	g.=	p.=
AAAATCCGTAGGGTTGC	2	25	ATGCTTTGTACGGGTGGT	g.=	p.=
AAAATTTGAAGTGGTAC	1	19	ATGCTTTGTACGGGTGGT	g.138G>A	p.=
AAAATTTGAAGTGGTAC	2	16	ATGCTTTGTACGGGTGGT	g.138G>A	p.=
AAAAGAACGCTCTAGTAT	1	96	ATGCTTTGTACGGGTGGT	g.23T>C	p.Leu8Pro
AAAAGAACGCTCTAGTAT	2	71	ATGCTTTGTACGGGTGGT	g.23T>C	p.Leu8Pro
AAAATAGAACCGTGGCAC	1	29	ATGCTTTGTACGGGTGGT	g.113C>T	p.Pro38Leu
AAAATAGAACCGTGGCAC	2	22	ATGCTTTGTACGGGTGGT	g.113C>T	p.Pro38Leu
AAACACTGGTTAGACCCAA	1	88	ATGCTTTGTACGGGTGGT	g.22C>G	p.Leu8Val



Text Preview 238 Lines

```

protein,count_time_1,count_time_2,score
p.=,23212,17456,0.7520248147509909
p.Ala25Glu,227,176,0.775330396475771
p.Ala25Pro,71,56,0.7887323943661971
p.Ala25Ser,42,33,0.7857142857142857
p.Ala25Thr,90,68,0.7555555555555555
p.Ala25Val,99,75,0.7575757575757576
p.Ala30Gly,3,2,0.6666666666666666
p.Ala30Pro,33,22,0.6666666666666666
p.Ala30Ser,75,53,0.7066666666666667
p.Ala30Thr,122,97,0.7950819672131147
p.Ala33Asp,18,12,0.6666666666666666
p.Ala33Gly,111,80,0.7207207207207207
  
```

## Example 3: FASTQ and VAMP-seq

*vampseq\_1\_1.fastq*

```
@COUNTESS-TEST-DATA:0
GGTAGAAGATAAAAAAAAGCA
+COUNTESS-TEST-DATA:0
ZZZZXYZZYXYZZYXYZZYX
@COUNTESS-TEST-DATA:1
GCTACTACGTGGACTGGCCA
+COUNTESS-TEST-DATA:1
ZZZZXYZZYXYZZYXYZZYX
@COUNTESS-TEST-DATA:2
AAGGGATCGAAAGTCCCAAT
+COUNTESS-TEST-DATA:2
ZZZZXYZZYXYZZYXYZZYX
```

*vampseq\_1\_2.fastq*

```
@COUNTESS-TEST-DATA:0
CTATATGGGCATTCCCGACT
+COUNTESS-TEST-DATA:0
ZZZZXYZZYXYZZYXYZZYX
@COUNTESS-TEST-DATA:1
TGGCAGATAACGGTCAAGTCA
+COUNTESS-TEST-DATA:1
ZZZZXYZZYXYZZYXYZZYX
@COUNTESS-TEST-DATA:2
GAGGACCTTGTCTTGAAAG
+COUNTESS-TEST-DATA:2
ZZZZXYZZYXYZZYXYZZYX
```

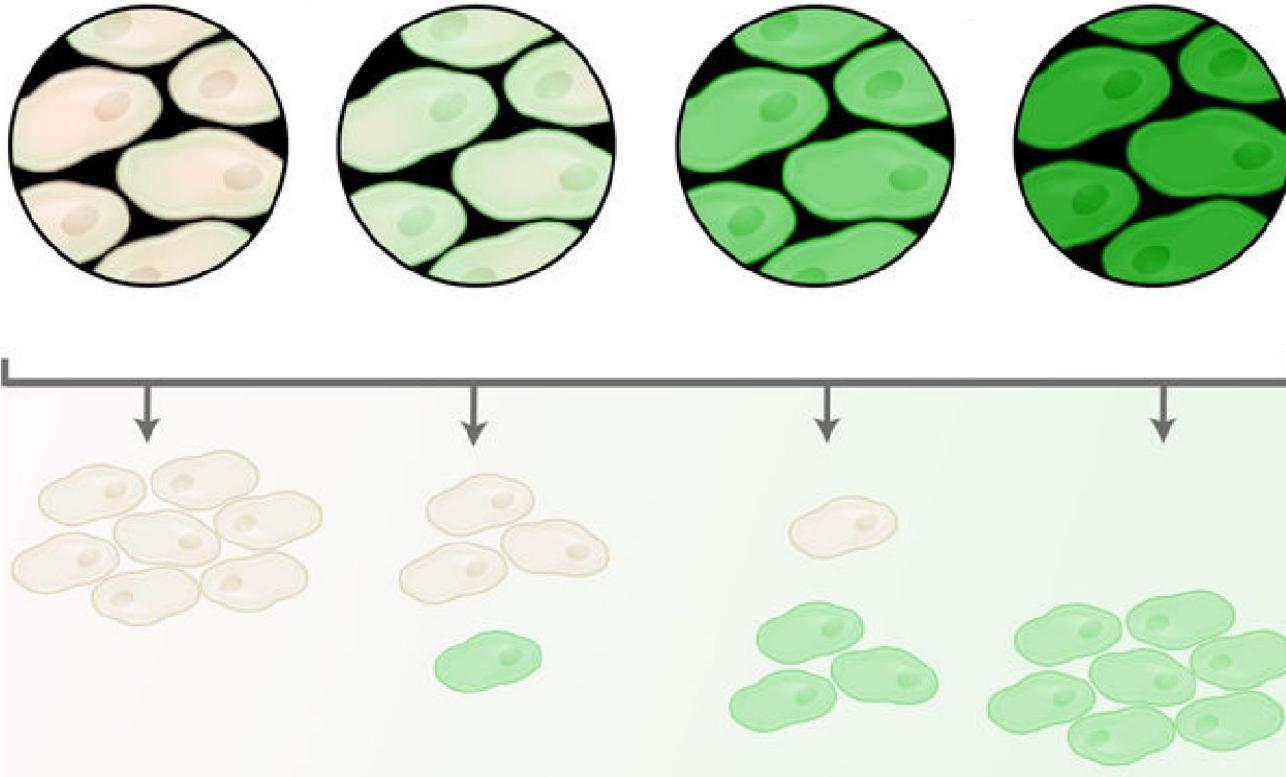
## Example 3: FASTQ and VAMP-seq

Cells are binned by GFP expression using FACS

Count variants in each bin by sequencing

Low abundance variant

WT-like abundance variant



Reference: ["Multiplex Assessment of Protein Variant Abundance by Massively Parallel Sequencing"](#)

### FASTQ Load

FASTQ Load 0.0.55 — Loads counts from FASTQ files containing either variant or barcodes

[add notes](#)

**Files** [+](#)

File	Filename	
File	vampseq_1_1.fastq	X
File	vampseq_1_2.fastq	X
File	vampseq_2_1.fastq	X
File	vampseq_2_2.fastq	X
File	vampseq_3_1.fastq	X
File	vampseq_3_2.fastq	X
File	vampseq_4_1.fastq	X
File	vampseq_4_2.fastq	X

Minimum Average Quality

Header Column?

Filename Column?

Group by Sequence?

Dataframe Preview 7563 rows

sequence object (index)	filename object (index)	count int64
AAAAATCCGTAGGGGTTGCC	vampseq_1_1.fastq	8
AAAAATCCGTAGGGGTTGCC	vampseq_1_2.fastq	9
AAAAATCCGTAGGGGTTGCC	vampseq_2_1.fastq	16
AAAAATCCGTAGGGGTTGCC	vampseq_2_2.fastq	17
AAAAATCCGTAGGGGTTGCC	vampseq_3_1.fastq	14
AAAAATCCGTAGGGGTTGCC	vampseq_3_2.fastq	22
AAAAATCCGTAGGGGTTGCC	vampseq_4_1.fastq	9
AAAAATCCGTAGGGGTTGCC	vampseq_4_2.fastq	9
AAAACTTGAAGTGGGTACG	vampseq_1_1.fastq	4

```

graph TD
    FASTQ[FASTQ Load] --> Regex[Regex Tool]
    Barcode[Load Barcode Map] --> Variant[Variant Translator]
    Variant --> Join((Join))
    Regex --> Join
    Join --> PivotBin[Pivot on bin]
    PivotBin --> CalcScore[Calculate Score]
    CalcScore --> PivotReplicate[Pivot on Replicate]
    PivotReplicate --> SaveScores[Save Scores]
  
```

## Regex Tool

Regex Tool 0.0.55 — Apply regular expressions to a column to make new column(s)

[add notes](#)

**Input Column:** filename

**Regular Expression:** vampseq\_(\d)\_(\d).fastq

**Output Columns:**

	Column Name	Column Type
Col	bin	integer
Col	rep	integer

**Drop Column:**

**Drop Unmatched Rows:**

**Multi Match:**

**Dataframe Preview 7563 rows**

sequence object (index) ▲	filename object (index) ▲	count int64 ▲	bin int64 ▲	rep int64 ▲
AAAAATCCGTAGGGGTTGCC	vampseq_1_1.fastq	8	1	1
AAAAATCCGTAGGGGTTGCC	vampseq_1_2.fastq	9	1	2
AAAAATCCGTAGGGGTTGCC	vampseq_2_1.fastq	16	2	1
AAAAATCCGTAGGGGTTGCC	vampseq_2_2.fastq	17	2	2
AAAAATCCGTAGGGGTTGCC	vampseq_3_1.fastq	14	3	1
AAAAATCCGTAGGGGTTGCC	vampseq_3_2.fastq	22	3	2
AAAAATCCGTAGGGGTTGCC	vampseq_4_1.fastq	9	4	1
AAAAATCCGTAGGGGTTGCC	vampseq_4_2.fastq	9	4	2
AAAATTGAACTGGGTACG	vampseq_1_1.fastq	4	1	1

```

graph TD
    A[FASTQ Load] --> C[Variant Translator]
    B[Load Barcode Map] --> C
    C --> D[Join]
    D --> E[Pivot on bin]
    E --> F[Calculate Score]
    F --> G[Pivot on Replicate]
    G --> H[Save Scores]
  
```

## Variant Translator

Variant Translator 0.0.55 — Turns a DNA sequence into a HGVS variant code

[add notes](#)

```

graph TD
    A[FASTQ Load] --> B[Regex Tool]
    C[Load Barcode Map] --> D[Variant Translator]
    B --> E[Join]
    D --> E
    E --> F[Pivot on bin]
    F --> G[Calculate Score]
    G --> H[Pivot on Replicate]
    H --> I[Save Scores]
    
```

**Input Column:** sequence

**Reference Column:** — NONE —

**\*OR\* Reference Sequence:** ATGCTTTGTACGGGTGGTGCCCTGGCTTATCTATCTAGATCGTCTCCGAGTCACGGTCAATTAGGTACTGCACTATCCTTGAGGC

**Output Column:** variant

**Max Mutations:** 10

**Protein Column:** protein

**Protein Offset:** 0

**Max Protein Variations:** 10

**Drop unidentified variants:**

**Drop Input Column(s):**

**Dataframe Preview 1000 rows**

barcode object	sequence object	variant object	protein object
CGTGTCTATTCGACCATCA	ATGCTTTGTACGGGTGGTGCCCTGGCTTA	g.100A>C	p.Thr34Pro
GGAGCATAGGTGGAGCTACA	ATGCTTTGTACGGGTGGTGCCCTGGCTTA	g.100A>C	p.Thr34Pro
CGAGTCGTCCTTCTCGTA	ATGCTTTGTACGGGTGGTGCCCTGGCTTA	g.100A>G	p.Thr34Ala
GCCAACAAACAGAACAGATGG	ATGCTTTGTACGGGTGGTGCCCTGGCTTA	g.100A>G	p.Thr34Ala
CATCATTTCATACTGAACCTC	ATGCTTTGTACGGGTGGTGCCCTGGCTTA	g.100A>G	p.Thr34Ala
CCCTGACCTCAATCTTCTG	ATGCTTTGTACGGGTGGTGCCCTGGCTTA	g.100A>T	p.Thr34Ser
AAGACGTAAGAGCTTGTAC	ATGCTTTGTACGGGTGGTGCCCTGGCTTA	g.101C>A	p.Thr34Lys
TGTACCCACCCACGCCGTCG	ATGCTTTGTACGGGTGGTGCCCTGGCTTA	g.101C>A	p.Thr34Lys
CTTAACCTCAGTTCTGCTAG	ATGCTTTGTACGGGTGGTGCCCTGGCTTA	g.101C>A	p.Thr34Lys

### Pivot on bin

Pivot Tool 0.0.55 — Groups a dataframe and pivots column values into columns.  
 Expanded column values are duplicated for each combination of pivot values. Missing values default to 0, and duplicate values are summed.

[add notes](#)

"sequence\_x"  
 Drop

"filename"  
 Drop

"count"  
 Expand

"bin"  
 Pivot

"rep"  
 Index

"sequence\_y"  
 Drop

"variant"  
 Drop

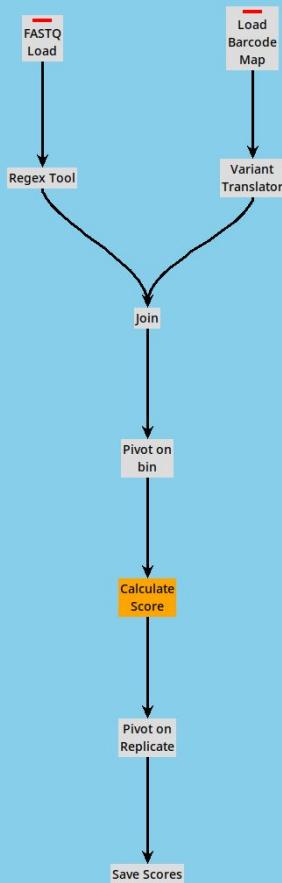
**Dataframe Preview 474 rows**

rep int64 (index)	protein object (index)	count_bin_1 float64	count_bin_2 float64	count_bin_3 float64	count_bin_4 float64
1	p.=	6099.	5783.	5641.	5068.
2	p.=	5936.	5898.	5537.	5145.
1	p.Ala25Glu	29.	40.	60.	66.
2	p.Ala25Glu	23.	42.	61.	73.
1	p.Ala25Pro	42.	15.	6.	2.
2	p.Ala25Pro	25.	17.	9.	0.
1	p.Ala25Ser	25.	22.	18.	42.
2	p.Ala25Ser	43.	23.	27.	37.

```

graph TD
    A[FASTQ Load] --> B[Regex Tool]
    C[Load Barcode Map] --> D[Variant Translator]
    B --> E[Join]
    D --> E
    E --> F[Pivot on bin]
    F --> G[Calculate Score]
    G --> H[Pivot on Replicate]
    H --> I[Save Scores]
  
```

## Calculate Score



Python Code 0.0.55 — Apply python code to each row.

Columns are mapped to local variables and back. If you assign to a variable called "`_filter`", only rows where that value is true will be kept.

[add notes](#)

```
score = (0.25 * count_bin_1 + 0.5 * count_bin_2 + 0.75 * count_bin_3 + count_bin_4) / (count_bin_1 + count_bin_2 + count_bin_3 + count_bin_4)
```

Python Code

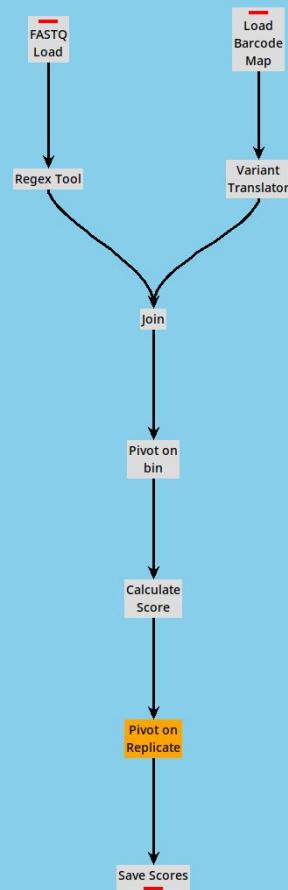
Drop Null Columns?



Dataframe Preview 474 rows

rep	protein	count_bin_1	count_bin_2	count_bin_3	count_bin_4	score
int64	object	float64	float64	float64	float64	float64
1	p.=	6099.	5783.	5641.	5068.	0.607100172635
1	p.Lys47Met	5.	9.	21.	19.	0.75
1	p.Lys47_Val48del	68.	49.	35.	15.	0.495508982036
1	p.Met1Arg	6.	11.	18.	11.	0.684782608696
1	p.Met1Ile	97.	74.	47.	58.	0.559782608696
1	p.Met1Leu	15.	40.	66.	94.	0.777906976744
1	p.Met1Lys	12.	25.	13.	1.	0.514705882353
1	p.Met1Thr	42.	39.	42.	20.	0.56993006993

## Pivot on Replicate



Pivot Tool 0.0.55 — Groups a dataframe and pivots column values into columns.

Expanded column values are duplicated for each combination of pivot values. Missing values default to 0, and duplicate values are summed.

[add notes](#)

"rep"

Pivot

"protein"

Index

"count\_bin\_1"

Drop

"count\_bin\_2"

Drop

"count\_bin\_3"

Drop

"count\_bin\_4"

Drop

"score"

Expand

Dataframe Preview 237 rows

protein object (index)	score_rep_1 float64	score_rep_2 float64
p.=	0.607100172635	0.609821904424
p.Ala25Glu	0.708974358974	0.731155778894
p.Ala25Pro	0.376923076923	0.421568627451
p.Ala25Ser	0.679906542056	0.611538461538
p.Ala25Thr	0.665816326531	0.618279569892
p.Ala25Val	0.545698924731	0.5225
p.Ala30Gly	0.34375	0.322916666667
p.Ala30Pro	0.34756097561	0.331818181818

# CountESS Status

Currently v0.0.61  
About to be v0.1.0

Working with several datasets and  
performing quite well.

Still needs feedback and I'm available  
to help!

# CountESS Future

Reading paired sequences.

Reading SRA datasets directly!

Simple visualizations.

More efficient processing for large  
data files.

?

