

## Summary of Notes:

- 1 – 3 : main databases (wrangling)
- 4 – 5 : merging and analyses
- 6 : graph
- 7: stats for graphs

## Summary for python code:

- 1 : Kew\_web\_scrape\_chromosome\_size.py
- 2 : ccdb\_edit.html
- 3 – 7: ccdb\_kew\_pw\_analyses.html

### 1.) KEW:

Web scrape to get Kew chromosome size:

Data taken from: <https://cvalues.science.kew.org/search>

Script: Kew\_web\_scrape\_chromosome\_size.py

File: Kew\_1C\_pg.csv

### 2.) CCDB: Chromosome Count DataBase

Originally taken from 3 files:

Angiosperms\_Jan2021.csv

Gymnosperms\_Jan2021.csv

Pterids\_Jan2021.csv

Cleaned and merge files:

- I. made sure to remove any unwanted whitespace in names
- II. created 3 columns:
  - A. major\_group: Angiosperms, Gymnosperms, Horsetails, Lycophytes, or Ferns
  - B. spory: heterosporous or homosporous
  - C. category: combine major\_group and spory group names for when we make graphs
- III. create a species column using the “resolved\_name” column
  - A. removed the 5 angiosperms with no resolved\_name from the database
  - B. nine where the species is a cross; I made the species nan
    - 1. (where ‘genus’ column startswith('x'))
  - C. 124 where ‘Ã’ is in the species name: change those to correct names
  - D. finally, there are still 19 Solanum genera with no species name (these happen to also have contained the ‘Ã’ within the resolved name – followed by author names. NO species)
  - E. Total of 28 rows of data where species is nan (9 cross and the 18 Solanum) in species column

Save this dataframe as:

**ccdb\_clean.csv**

377,563 rows x 20 columns

see: [ccdb\\_edit.html](#)

Made a mini version of this with 10 columns:

**ccdb\_clean\_mini.csv**

10 cols are: ['family', 'genus', 'species', 'resolved\_name', 'major\_group', 'spory', 'category', 'parsed\_n', 'list\_parsed\_n', 'min\_parsed\_n']

(NOTE: parsed\_n, and subsequently the other n columns, have 10,035 rows of missing data)

### 3.) PW's 10 sample database

Data from this table will be added to kew where there's info in the 1C column, and to ccdb where there's info in the n column.

**"Paul\_chromosome\_1C\_additions.csv"**

(Original file that I didn't use is: **Paul\_orig\_chromosome\_1C\_additions.csv**)

Citations:

Heterosporous ferns:

Li, F.-W., Brouwer, P., Carretero-Paulet, L., Cheng, S., de Vries, J., Delaux, P.-M., et al. (2018). Fern genomes elucidate land plant evolution and cyanobacterial symbioses. *Nat Plants* 4, 460–472.

Sellaginella:

Little DP, Moran RC, Brenner ED, Stevenson DW. Nuclear genome size in Selaginella. *Genome*. 2007 Apr;50(4):351–6. doi: 10.1139/g06-138. PMID: 17546093.

genus	sp	category	n	DNA Amount1C (pg)
<b>Azolla</b>	filicoides	Ferns\nHeterosporous		0.76
<b>Salvinia</b>	cucullata	Ferns\nHeterosporous	9	0.26
<b>Pilularia</b>	americana	Ferns\nHeterosporous		0.84
<b>Regnellidium</b>	diphyllum	Ferns\nHeterosporous	19	1.8
<b>Marsilea</b>	spp	Ferns\nHeterosporous	20	1.37
<b>Selaginella</b>	apoda	Lycophytes\nHeterosporous	9	
<b>Selaginella</b>	invovens	Lycophytes\nHeterosporous	9	
<b>Selaginella</b>	moellendorffii	Lycophytes\nHeterosporous	9	
<b>Selaginella</b>	vogelii	Lycophytes\nHeterosporous	9	
<b>Selaginella</b>	willdenowii	Lycophytes\nHeterosporous	9	

## 4.) Add pw data to ccdb and kew databases

merge kew + pw = **Kew\_pw\_1C\_pg\_mini.csv**

- columns: ['genus', 'species', 'DNA Amount1C (pg)']
- size: (12278, 3)

merge cccb (no horsetails, no missing n info (min\_parse\_n is nan)) + pw =  
**ccdb\_pw\_noHorsetails\_noMissingN.csv**

- columns: ['family', 'genus', 'species', 'resolved\_name', 'major\_group', 'spory', 'category', 'parsed\_n', 'list\_parsed\_n', 'min\_parsed\_n']
- size: (367405, 10)

## 5.) Databases for the Graphs

NOTE: This gets used to create the upper two graphs

Using the cccb\_pw\_noHorsetails\_noMissingN.csv, get only ONE minimum per genus:  
(final\_ccdb\_mins in html file)

**ccdb\_noHorsetails\_1min\_per\_genus.csv**

size: (8338, 10)

columns: ['genus', 'min\_parsed\_n', 'family', 'species', 'resolved\_name', 'major\_group', 'spory', 'category', 'parsed\_n', 'list\_parsed\_n']

NOTE: Did NOT end up using this one for making the graph.

Using the cccb\_pw\_noHorsetails\_noMissingN.csv, get ALL possible minimums per genus:

**ccdb\_noHorsetails\_ALLmin\_per\_genus.csv**size:

size: (69098, 4)

columns: ['genus', 'species', 'category', 'min\_parsed\_n']

NOTE: This does NOT used for making the lower 2 graphs

Want all diploids, even if aneuploid....up to a certain point:

Want to groupby each genus. Within each genus, make a list of all chromosome counts. Within the list of counts, divide each number in list by the lowest value in the list. If that divided value is  $\leq 1.2$ , then we will keep those individuals within that given genus.

This contains all unique species for each genus!

**ccdb\_cleaned EVERY\_min\_1pt2\_ratio.csv**

shape: (98373, 11)

```
columns: ['family', 'genus', 'species', 'resolved_name', 'major_group',
'spory', 'category', 'parsed_n', 'list_parsed_n', 'min_parsed_n',
'diploid_mins']
```

NOTE: This gets used for the lower graphs

Use ccdb\_cleaned EVERY\_min\_1pt2\_ratio.csv.

Merge with ccdb\_pw\_noHorsetails\_noMissingN.csv (kew\_pw2 from html).

Select just needed columns: ['genus', 'species', 'category', 'min\_parsed\_n', 'DNA Amount1C (pg)']

Drop the duplicate lines.

Result:

**lower\_graphs\_db.csv**

(3656, 5)

## 6.) Graphs:

upper graphs use: **ccdb\_noHorsetails\_1min\_per\_genus.csv**

lower graphs use: **lower\_graphs\_db.csv**

output graph:

**hetero\_chromo\_num\_paper\_withPWdata\_ALLmins\_ratio\_1pt2.png**

## 7.) Stats:

### UPPER GRAPHS: chromosome counts

CATEGORY

Angiosperms\nHeterosporous 7900

Ferns\nHomosporous 343

Gymnosperms\nHeterosporous 76

Lycophytes\nHomosporous 7

Lycophytes\nHeterosporous 6

Ferns\nHeterosporous 6

Name: category, dtype: int64

HOMOSPOROUS ONLY

count 350.000000

mean 40.471429

std 23.980365

min 1.000000

25% 30.000000

50% 36.000000  
75% 41.000000  
max 164.000000  
Name: min\_parsed\_n, dtype: float64

HETEROSPOROUS ONLY  
count 7988.000000  
mean 12.878443  
std 8.842195  
min 1.000000  
25% 9.000000  
50% 11.000000  
75% 15.000000  
max 275.000000  
Name: min\_parsed\_n, dtype: float64

BOTH (HETEROSPOROUS & HOMOSPOROUS)  
count 8338.000000  
mean 14.036699  
std 11.384060  
min 1.000000  
25% 9.000000  
50% 11.000000  
75% 16.000000  
max 275.000000  
Name: min\_parsed\_n, dtype: float64

## LOWER GRAPHS: genome size

CATEGORY  
Angiosperms\nHeterosporous 3387  
Gymnosperms\nHeterosporous 178  
Ferns\nHomosporous 72  
Lycophytes\nHeterosporous 14  
Ferns\nHeterosporous 4  
Lycophytes\nHomosporous 1  
Name: category, dtype: int64

HOMOSPOROUS ONLY  
count 73.000000  
mean 18.356027

std 21.938207  
min 2.430000  
25% 6.970000  
50% 10.460000  
75% 14.890000  
max 74.840000  
Name: DNA Amount1C (pg), dtype: float64

#### HETEROSPOROUS ONLY

count 3583.000000  
mean 5.730851  
std 9.370177  
min 0.080000  
25% 0.750000  
50% 1.750000  
75% 6.165000  
max 102.900000  
Name: DNA Amount1C (pg), dtype: float64

#### BOTH (HETEROSPOROUS & HOMOSPOROUS)

count 3656.000000  
mean 5.982940  
std 9.932141  
min 0.080000  
25% 0.770000  
50% 1.815000  
75% 6.700000  
max 102.900000  
Name: DNA Amount1C (pg), dtype: float64