



How We Manage Large-Scale Data Collection

Darryl Chamberlain, Ph.D.

Emily Faulconer, Ph.D.

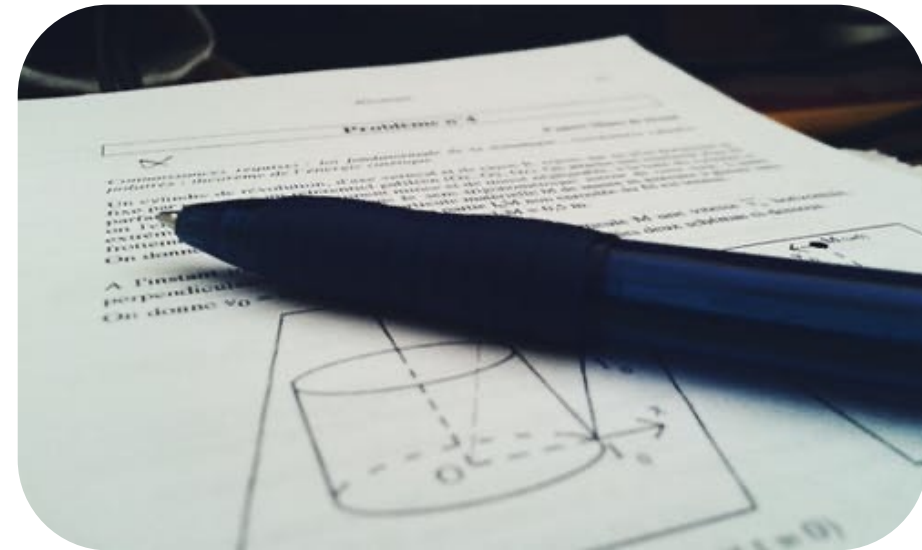
Embry-Riddle Aeronautical University

As a mixed-method study, we are collecting a variety of data across 18 terms.

All data is collected in each term of MATH 111 and PHYS 102 between August 2021 and September 2023.

- Withdrawals in MATH 111 and PHYS 102
- Responses to voluntary, 8 Likert scale question survey
- Student performance data (discussion grades and final grade)
- Transcripts of discussion responses
- Focus-group interviews

For context, PHYS 102 AUG '21 had 182 students and 39,682 sentences in discussion responses.



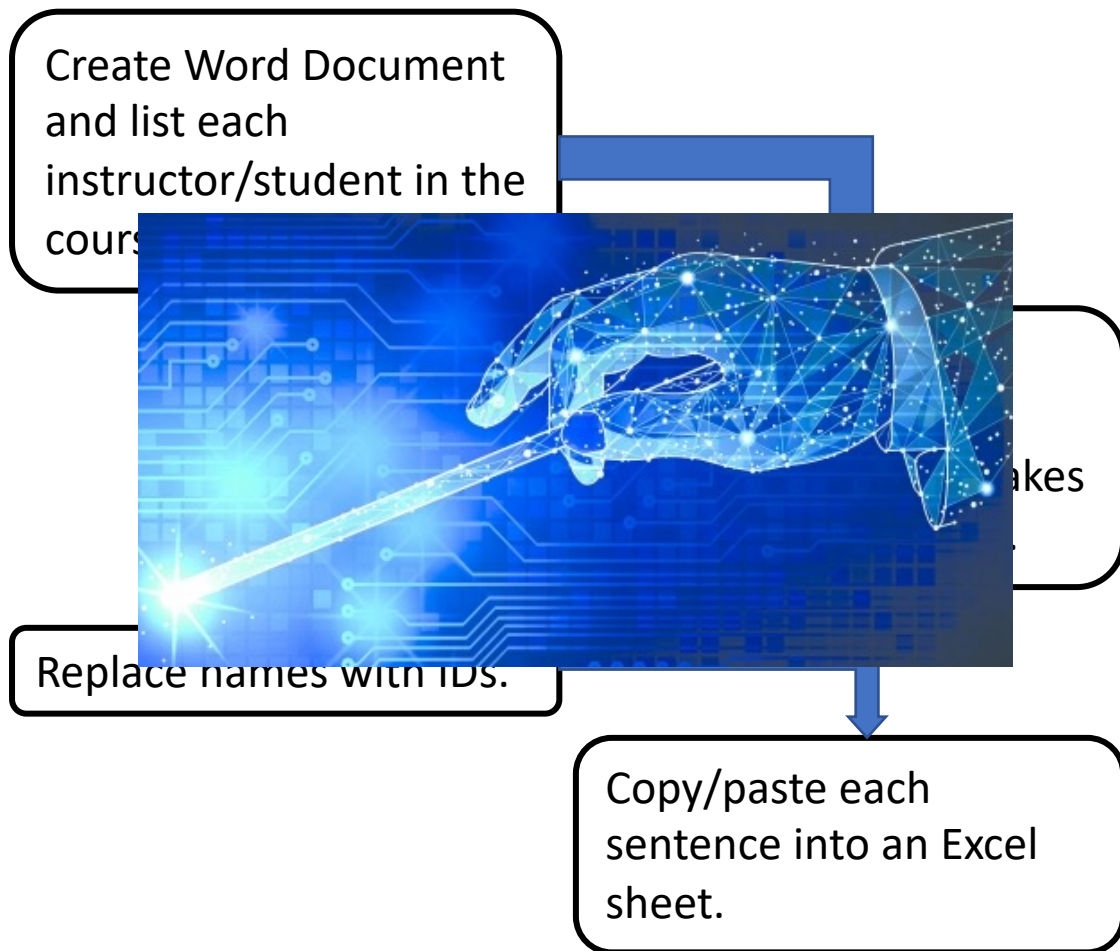
Data cleaning involves standardizing data for analysis and removing incomplete, irrelevant, or erroneous data.

For our data, this primarily included:

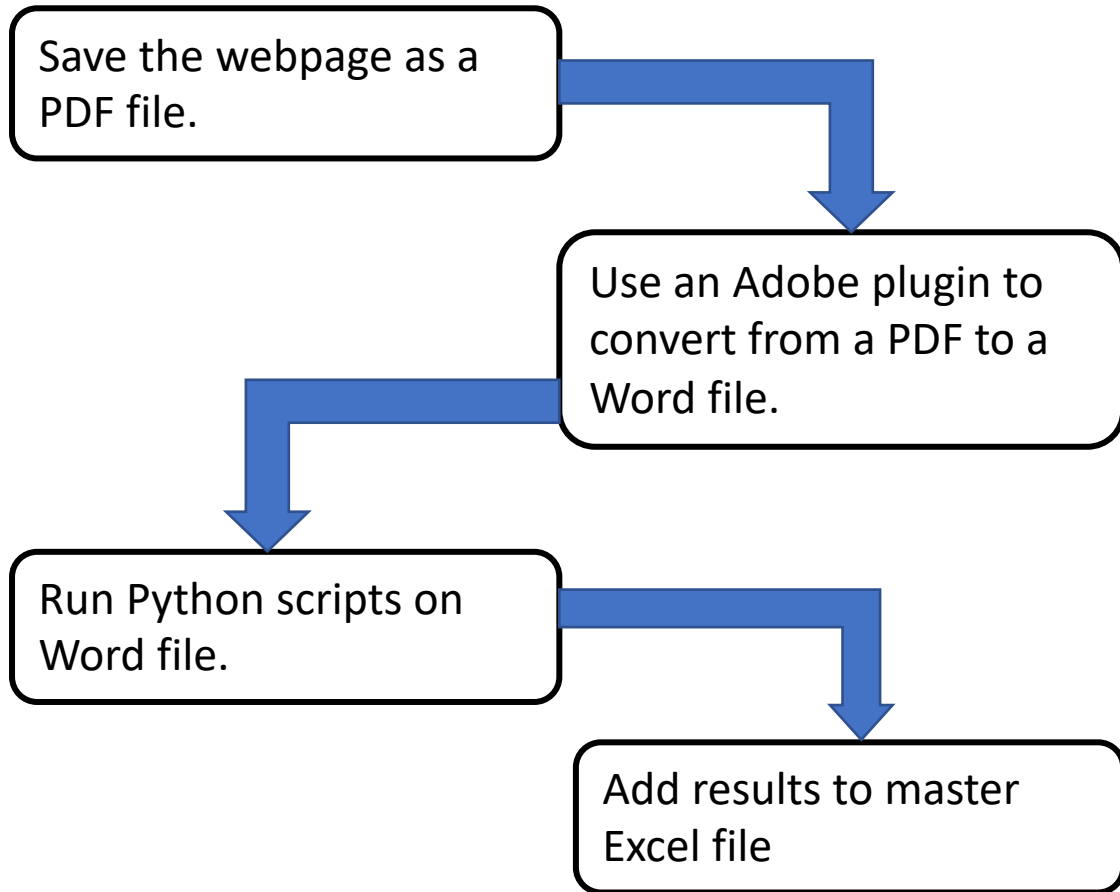
- Standardizing withdrawal/drop reports
- Downloading discussion transcripts and parsing them by sentence for analysis
- Organizing all data for cross-semester analysis



Initially, we manually collected and cleaned discussion transcripts.



Now, we let technology do the heavy lifting.



A simple trick to improve your data organization is to put all data into columns.

	A	B	C	D	E	F	H
1	Speaker	Speaker Type	Class ID	Term	Cohort	Module Number	Sentence
13583	QVXCAP	Student	PHY	AUG21	2	8	In the principles of electromagnetism, a magnetic field is induced by a changing electric field.
13584	QVXCAP	Student	PHY	AUG21	2	8	When a magnet moves through a coil of wires, electric current is induced as long as the magnet is moving.
13585	QVXCAP	Student	PHY	AUG21	2	8	Magnetic flux, a measurement of a magnetic field in a given area, depends on the size of the area.

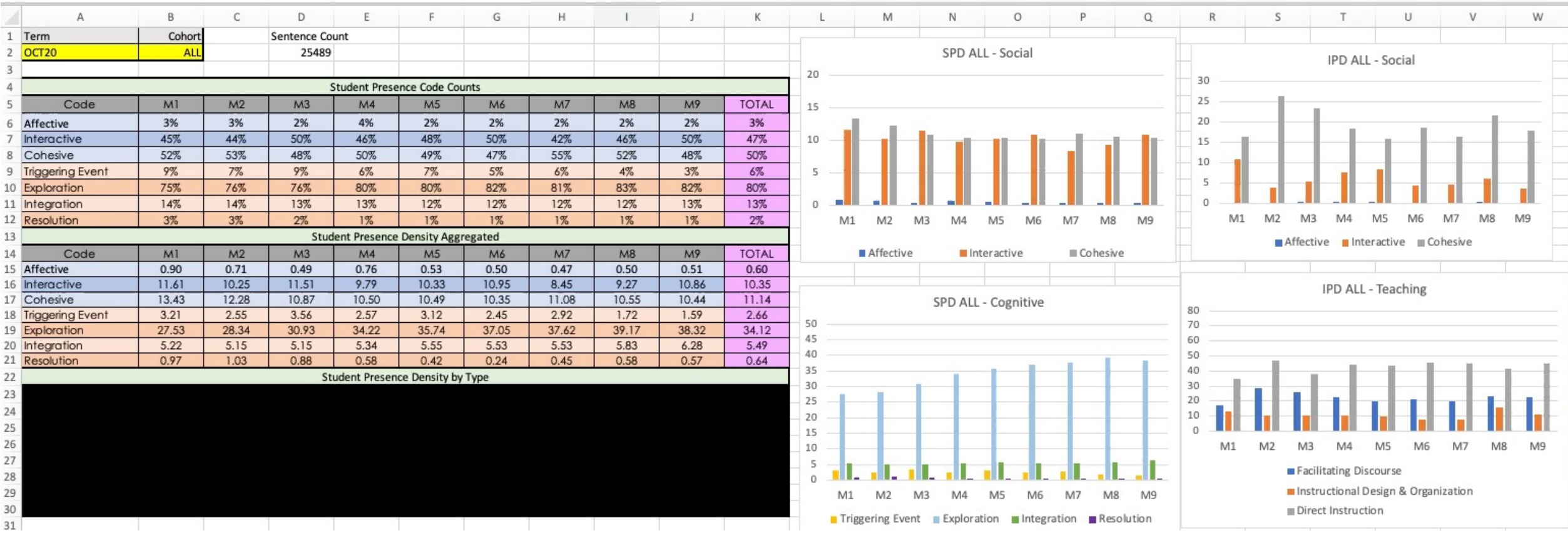
I	J	K	L	M	N	O	P	Q	S
Word Count	Emily	Syaza	Abigail	Patricia	Amina	Qaisar	Andre	Negotiated Code	COMMENTS
27				IS			IS		
34				IS			IS		
39				IS			IS		

Common calculations such as AVERAGE, STDEV, and COUNTIF are far easier to complete when you can call an entire column.

```
=IF(AND(OR($C$1="Emily", $C$2="Emily"), OR($C$1="Syaza", $C$2="Syaza")), IF($C$3="ALL", CO  
IF(AND(OR($C$1="Abigail", $C$2="Abigail"), OR($C$1="Syaza", $C$2="Syaza")), IF($C$3="ALL", C  
IF(AND(OR($C$1="Amina", $C$2="Amina"), OR($C$1="Qaisara", $C$2="Qaisara")), IF($C$3="ALL  
IF(AND(OR($C$1="Patrick", $C$2="Patrick"), OR($C$1="Andrea", $C$2="Andrea")), IF($C$3="ALL
```

```
COUNTIFS(DATA!$J:$J, IRR!$B22, DATA!$K:$K, IRR!$N$5), COUNTIFS(DATA!$D:$D, $C$3, DATA!$J:$J, IRR!$B22, DATA!$K:$K, IRR!$N$5)),  
", COUNTIFS(DATA!$L:$L, IRR!$B22, DATA!$K:$K, IRR!$N$5), COUNTIFS(DATA!$D:$D, $C$3, DATA!$L:$L, IRR!$B22, DATA!$K:$K, IRR!$N$5)),  
ALL", COUNTIFS(DATA!$N:$N, IRR!$B22, DATA!$O:$O, IRR!$N$5), COUNTIFS(DATA!$D:$D, $C$3, DATA!$N:$N, IRR!$B22, DATA!$O:$O, IRR!$N$5)),  
ALL", COUNTIFS(DATA!$M:$M, IRR!$B22, DATA!$P:$P, IRR!$N$5), COUNTIFS(DATA!$D:$D, $C$3, DATA!$M:$M, IRR!$B22, DATA!$P:$P, IRR!$N$5)), "NA"))))
```

Strong data organization can allow seamless data analysis.



Data cleaning has a direct impact on data analysis.

[illegible]