

Cleaning the Ruhlman Files

Megan O, 5/3/16

I played a part in converting the original Word files to dirty CSV files. These CSVs (one for each year) were the starting point files that others used to tackle the more fine-tuned cleaning, given that each document was formatted slightly differently and sometimes had different information. See `parser.py`.

I was then tasked to clean the 2012 dirty CSV. I ended up doing this manually in Google Sheets (see README for the spreadsheet link)- 2012 wasn't too bad, and most of the cleaning consisted of moving presenter names out of the abstracts of the previous project (due to a glitch in the parser) and to the correct project.

Finally, I wrote a few scripts (before all the years got cleaned separately) to clean up the abstracts so that I could generate that text file with all the abstracts. See `abstracts.py` and `absRevised.py`.