the old one was something like this, where we followed the links in the 2nd and 3rd column, remember?

http://www.governotransparente.com.br/transparencia/4382489/consultarliqdesporc/resultado?
ano=7&inicio=15%2F09%2F2020&fim=15%2F09%2F2020&orgao=-1&elem=-1&unid=-1&valormax=&valormin=&credor=-
1&clean=false&datainfo=MTIwMjEwMjI4MTk1OVBQUA%3D%3D

Well now the links are not there anymore and we need to use this page which has the 2nd and 3rd columns merged for some
reason, showing two things in one cell. And the links are in the last column.

https://www.governotransparente.com.br/transparencia/4409486/consultarliqdesporc?
ano=5&inicio=01%2F01%2F2021&fim=22%2F03%2F2021&unid=-1&valormax=&valormin=&orgao=-1&elem=-1&credor=-
1&clean=false&datainfo=MTIwMjEwMzIyMTYxN1BQUA%3D%3D

What I need is to be able to scrape the same information. This means
- using the new links (icons)
- scrape the 4 things in the column 2 and 3: Empenho, Unidade gestora, Documento, Natureza da despesa
- scrape the rest as usual, based on the exisiting scraper

I have a new version of the Scraper that I modified to my needs so I will send it to you so you can start with it

http://www.governotransparente.com.br/acessoinfo/4382489/consultarcontratoaditivo/resultado/detalhes/4651/0?clean=false

https://www.governotransparente.com.br/acessoinfo/4382489/consultarcontratoaditivo/resultado/detalhes/1919/1?
clean=false

you can see that these two pages have links in the data.
One of them is in a table, in a column, so that's easier to organize. We take the column name and put scraped text there,
then we can take the column name with the suffix _link to put the links

the other is probably going to be the same thing. We can call the columns Text and Link

I also plan to split the scraper into two. You can notice that we first scrape some main table and then we go into the details
and scrape the details.
I am planning on having these two scrapers separate and independent. The first scraper is already writing the links-to-be-
scraped into the file, so the second scraper can read them and scrape.

On a local machine, it is sufficient if the scripts are separate, we initialize the first script, get the scraped data that include the
links.
Then we initialize the second scraper that goes through the links and scrapes its part.

actually, I can see that we don't store the links in the main table files. Just in the detailed files.
Can we also save the links in the main table files? It also makes sense since now there are icons instead of the text I mean
the links that lead to the detailed pages

rephrasing again, to be clear....
as of now, we just pass the links in an array to the second part of the scraper and then the second part of the scraper uses
them to get the website and writes these links in the detailed files as a second column. I would like these links also appearing
in the main table

http://www.governotransparente.com.br/acessoinfo/4382489/consultarcontratoaditivo/resultado/detalhes/4651/0?clean=false

https://www.governotransparente.com.br/acessoinfo/4382489/consultarcontratoaditivo/resultado/detalhes/1919/1?
clean=false

as to these two links that I shared earlier. This is what I currently also need to scrape. The Main table is this one
http://www.governotransparente.com.br/acessoinfo/4382489/consultarcontratoaditivo/resultado?
ano=8&inicio=01%2F01%2F2021&fim=24%2F03%2F2021&contr=&credor=-1&clean=false

and it is scraped correctly. But the inner links are not well scraped. I know they were not part of the requirements last time. Can we include them now? There is some data with the link. One of the problems is that the current scraper does not recognize the sections correctly.
You can try it yourself, running this
python Scrape.py "consultarcontratoaditivo" "http://www.governotransparente.com.br/acessoinfo/4382489/consultarcontratoaditivo/resultado?ano=8&inicio=01%2F03%2F2021&fim=31%2F03%2F2021&contr=&credor=-1&clean=false" "2021-03" --dir "out_contratos_2021_03_test" --con 4

I included h5 into the scraper, but still wasn't able to solve it


as to the basic scraper, this is the format to run it. It creates the out folder and re-uses it if it exists

python Scrape.py "consultarliqdesporc" "http://www.governotransparente.com.br/transparencia/4382489/consultarliqdesporc?ano=8&inicio=01%2F01%2F2021&fim=25%2F03%2F2021&unid=-1&valormax=&valormin=&orgao=-1&elem=-1&credor=-1&clean=false&datainfo=MTIwMjEwMzI1MDM0MFBQUA%3D%3D" "2021-010203" --dir "out_2021_run1" --con 4


the current scraper also writes a log file in the out folder, in case we need to do any debugging

I would like to make the scraper work for both the new site format with the icons and the old format (where the links are in the second column with the number)