

## Website Input - Task #1809

### Refactor scrape\_page call

04/04/2017 05:24 PM - Luke Murphey

<b>Status:</b>	Closed	<b>Start date:</b>	04/04/2017
<b>Priority:</b>	Normal	<b>Due date:</b>	
<b>Assignee:</b>	Luke Murphey	<b>% Done:</b>	100%
<b>Category:</b>	Input: Web Spider	<b>Estimated time:</b>	0.00 hour
<b>Target version:</b>	4.2		
<b>Description</b>			
The scrape_page call takes a ton of arguments which makes it unwieldy.			

#### Associated revisions

---

**Revision 398 - 04/04/2017 06:26 PM - lukemurphey**

Separating scraping from the input code

Reference #1809

**Revision 401 - 04/04/2017 08:09 PM - lukemurphey**

Moving parameters to constructor & function calls

Reference #1809

**Revision 407 - 04/04/2017 08:09 PM - lukemurphey**

Moving parameters to constructor & function calls

Reference #1809

**Revision 408 - 04/05/2017 06:10 AM - lukemurphey**

Refactored controller and input

Redesigning the controller and input to work with the refactored code

Reference #1809

#### History

---

**#1 - 04/04/2017 05:25 PM - Luke Murphey**

Some options:

1. Make a separate scraper class
2. Pass arguments as a config object
3. Pass arguments as \*args or \*\*kwargs

**#2 - 04/04/2017 05:34 PM - Luke Murphey**

Here are the functions invoked by scrape\_page():

- resolve\_proxy\_type
- get\_result\_single
  - get\_result\_built\_in\_client
    - detect\_encoding
  - get\_result\_browser
    - get\_display
    - get\_firefox\_profile
    - add\_auth\_to\_url
  - unescape
  - escape\_field\_name
  - extract\_links
    - cleanup\_link

### #3 - 04/04/2017 05:37 PM - Luke Murphey

Here are the functions that are modular input specific:

- get\_file\_path
- run

I think the best approach is to

1. break out the modular input and leave everything else in the scraper
2. make the arguments to scrape\_page class variables
3. updated the controller, search command, tests

### #4 - 04/04/2017 05:41 PM - Luke Murphey

What should be in the constructor call versus the command call?

- **basics:** url, username, password, etc. *(in scrape\_page call)*
- **output customization:** output\_matches\_as\_mv, output\_matches\_as\_separate\_fields, etc. *(in constructor)*
- **charset detection:** charset\_detect\_meta\_enabled, etc.
- **proxy info:** proxy\_type, etc.
- **spider info:** page\_limit, etc. *(in constructor)*

### #5 - 04/04/2017 08:08 PM - Luke Murphey

Phases:

1. Break out WebScraper class
2. Add constructor options for proxy info and charset
  1. Modify get\_result\_browser, get\_result\_single, get\_result\_built\_in\_client, get\_http\_client
3. Update search command

4. Update modular input
5. Update controller

**#6 - 04/05/2017 06:48 AM - Luke Murphey**

- *Status changed from New to Closed*
- *Target version set to 4.2*
- *% Done changed from 0 to 100*