

# DOCX Standard Operating Procedure (SOP) Parser for SFB1208 Experiments

---

This repository contains a set of files to parse SOP from lab experiments.

## Why?

---

As research lab usually has their own set of SOP to conduct experiments, a tool to extract metadata from an editable document (e.g., DOCX) would be handy. The metadata is helpful in documenting the research and hence improves the reproducibility of the conducted research. To enable the metadata extraction, the SOP should follow some annotation rules (described later below).

## How is this repo structured?

---

- The base directory contains the metadata extraction script.
- *input* directory contains the docx SOPs example for extraction.
- *output* directory contains extracted steps order – key – value in both JSON and XLSX format.

## What is extracted from the SOP, and how is it represented in the docx document?

---

The parser extracts:

Extracted Items	Description	Representation	Example	Extracted order,key,value
Section	The section name	<section   section name>	<section   Structure Preparation>	<ul style="list-style-type: none"> <li>"-", section, Structure Preparation</li> </ul>
Order	The <i>order</i> of the steps, based on the order of the paragraph in the docx SOP document	-	-	-
Key	The <i>key</i> for the metadata, based on the value represented in the curly bracket after the pipe character {value   key}.	{value   key}	{sequence alignment   stage}	<ul style="list-style-type: none"> <li>&lt;order&gt;, stage, sequence alignment</li> </ul>
Comment	<i>Comments</i> are allowed within the key, represented within regular brackets after the pipe symbol.	{value   (comment) key}	{receptor residue   (minimization) target}	<ul style="list-style-type: none"> <li>&lt;order&gt;, (minimization) target, receptor residue</li> </ul>
Value	The <i>value</i> of the metadata is based on the first value represented in the curly bracket before the pipe character {value   key}. Example: with 'sequence alignment' as the value.	{value   key}	{sequence alignment   stage}	<ul style="list-style-type: none"> <li>&lt;order&gt;, stage, sequence alignment</li> </ul>
Control flow: <i>for each</i>	Extract multiple key value pairs related to <i>for each</i> iterations	<flow type   iterated value>	<for each   generated pose>	<ul style="list-style-type: none"> <li>&lt;order&gt;, step type, iteration</li> <li>&lt;order&gt;, flow type, for each</li> <li>&lt;order&gt;, flow parameter, generated pose</li> </ul>

Extracted Items	Description	Representation	Example	Extracted order,key,value
Control flow: <i>while</i>	Extract multiple key value pairs related to <i>while</i> iteration	<flow type key logical operator value> ... <increment/decrement operation increment/decrement value>	<while pH lte 7> ... <+ 1>	<ul style="list-style-type: none"> <li>• &lt;order&gt;, step type, <i>iteration</i></li> <li>• &lt;order&gt;, flow type, <i>while</i></li> <li>• &lt;order&gt;, flow parameter, pH</li> <li>• &lt;order&gt;, flow logical parameter, lte</li> <li>• &lt;order&gt;, flow compared value, 7</li> <li>• &lt;order&gt;, flow operation, +,</li> <li>• &lt;order&gt;, flow magnitude, 1</li> </ul>
Control flow: <i>if</i>	Extract multiple key value pairs related to <i>if</i> iteration	<if key logical operator value>	<if pH lte 7>	<ul style="list-style-type: none"> <li>• &lt;order&gt;, step type, <i>conditional</i></li> <li>• &lt;order&gt;, flow type, <i>if</i></li> <li>• &lt;order&gt;, flow parameter, pH</li> <li>• &lt;order&gt;, flow logical parameter, lte</li> <li>• &lt;order&gt;, flow compared value, 7</li> </ul>

Extracted Items	Description	Representation	Example	Extracted order,key,value
Control flow: <i>else if</i>	Extract multiple key value pairs related to <i>else if</i> iteration	<else if key logical operator value>	<else if pH between [8-12]>	<ul style="list-style-type: none"><li>• &lt;order&gt;, step type, <i>conditional</i></li><li>• &lt;order&gt;, flow type, <i>else if</i></li><li>• &lt;order&gt;, flow parameter, pH</li><li>• &lt;order&gt;, flow logical parameter, between</li><li>• &lt;order&gt;, flow range, [8-12]</li><li>• &lt;order&gt;,start iteration value,8</li><li>• &lt;order&gt;,end iteration value,12</li></ul>
Control flow: <i>else</i>	Extract multiple key value pairs related to <i>else</i> iteration	<else>		<ul style="list-style-type: none"><li>• &lt;order&gt;, step type, <i>conditional</i></li><li>• &lt;order&gt;, flow type, <i>else</i></li></ul>

Extracted Items	Description	Representation	Example	Extracted order,key,value
Control flow: <i>for</i>	Extract multiple key value pairs related to <i>for</i> iteration	<for   key   [range]   iteration_operation   magnitude>	<for   pH   [1-7]   +   1>	<ul style="list-style-type: none"> <li>• &lt;order&gt;, step type, <i>iteration</i></li> <li>• &lt;order&gt;, flow type, <i>for</i></li> <li>• &lt;order&gt;, flow parameter, pH</li> <li>• &lt;order&gt;, flow logical parameter, lte</li> <li>• &lt;order&gt;, flow flow range, [1-7]</li> <li>• &lt;order&gt;,start iteration value,1</li> <li>• &lt;order&gt;,end iteration value,7</li> <li>• &lt;order&gt;, flow operation, +,</li> <li>• &lt;order&gt;, flow magnitude, 1</li> </ul>

The overall example of the SOP document is available in the *input/sop2.docx* file. The color in the *sop2.docx* does not play any role in the order/key/value extraction.

## How the parser should be run?

1. Create SOP according to the above annotation rules.
2. Change the input directory/file name in the python script (2nd last line).
3. Change the output directory/filename (last line).

4. Run the script.

## What are the further plans?

---

1. Fixes for while control flow, and logical operators in general control flow.
2. Consult CAi and Biochemistry1 for its implementability on other labs.
3. Align the used keys with terms from an ontology, or if the term does not exist, create a new term by extending an ontology or creating a term within a new ontology.