

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 <i>Structure Preparation</i> > | <ul style="list-style-type: none">• "-", section, <i>Structure Preparation</i> |
| 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 <i>stage</i> } | <ul style="list-style-type: none">• <order>, <i>stage</i>, sequence alignment |
| Comment | <i>Comments</i> are allowed within the key, represented within regular brackets after the pipe symbol. TBD : How to serialize comments in the metadata parser output. | {value (comment) key} | {receptor residue <i>(minimization)</i> target} | <ul style="list-style-type: none">• <order> target, receptor residue |

| EXTRACTED ITEMS | DESCRIPTION | REPRESENTATION | EXAMPLE | EXTRACTED ORDER,KEY,VALUE |
|----------------------------------|---|--|-------------------------------|--|
| 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">• <order>, stage, sequence alignment |
| 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">• <order>, step type, iteration• <order>, flow type, for each• <order>, flow parameter, generated pose |
| 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">• <order>, step type, iteration• <order>, flow type, while• <order>, flow |

| EXTRACTED ITEMS | DESCRIPTION | REPRESENTATION | EXAMPLE | EXTRACTED ORDER,KEY,VALUE |
|--------------------|-------------|----------------|---------|---|
| | | | | parameter, pH |
| | | | | • <order>, flow logical parameter, lte |
| | | | | • <order>, flow compared value, 7 |
| | | | | • <order>, flow operation, +, |
| | | | | • <order>, flow magnitude, 1 |

| EXTRACTED ITEMS | DESCRIPTION | REPRESENTATION | EXAMPLE | EXTRACTED ORDER,KEY,VALUE |
|----------------------------|---|---------------------------------|---------------|---|
| 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">• <order>, step type, <i>conditional</i>• <order>, flow type, <i>if</i>• <order>, flow parameter, pH• <order>, flow logical parameter, lte• <order>, flow compared value, 7 |

| 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">• <order>, step type, <i>conditional</i>• <order>, flow type, <i>else if</i>• <order>, flow parameter, pH• <order>, flow logical parameter, between• <order>, flow range, [8-12]• <order>,start iteration value,8• <order>,end iteration value,12 |

| EXTRACTED ITEMS | DESCRIPTION | REPRESENTATION | EXAMPLE | EXTRACTED ORDER,KEY,VALUE |
|------------------------------|---|---|--------------------|--|
| Control flow: <i>else</i> | Extract multiple key value pairs related to <i>else</i> iteration | <else> | | <ul style="list-style-type: none">• <order>, step type, <i>conditional</i>• <order>, flow type, <i>else</i> |
| 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">• <order>, step type, <i>iteration</i>• <order>, flow type, <i>for</i>• <order>, flow parameter, pH• <order>, flow logical parameter, lte• <order>, flow flow range, [1-7] |

| EXTRACTED ITEMS | DESCRIPTION | REPRESENTATION | EXAMPLE | EXTRACTED ORDER,KEY,VALUE |
|--------------------|-------------|----------------|---------|--|
| | | | | <ul style="list-style-type: none">• <order>,start iteration value,1• <order>,end iteration value,7• <order>, flow operation, +,• <order>, flow magnitude, 1 |

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.