

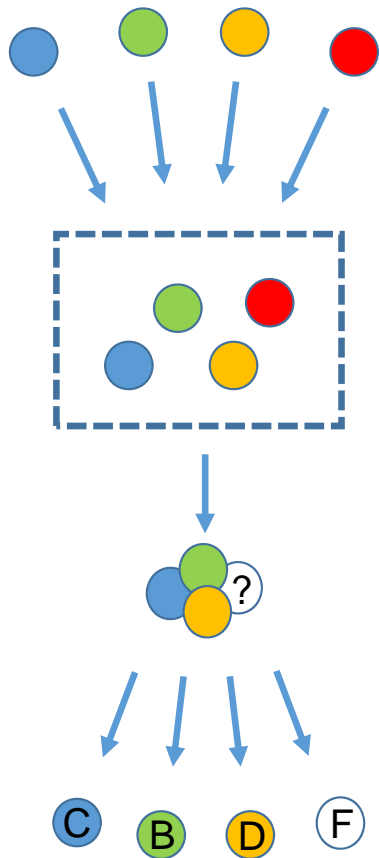
# Protein-Protein Interactions

## *Introduction to Interaction Data Processing*

February 2020

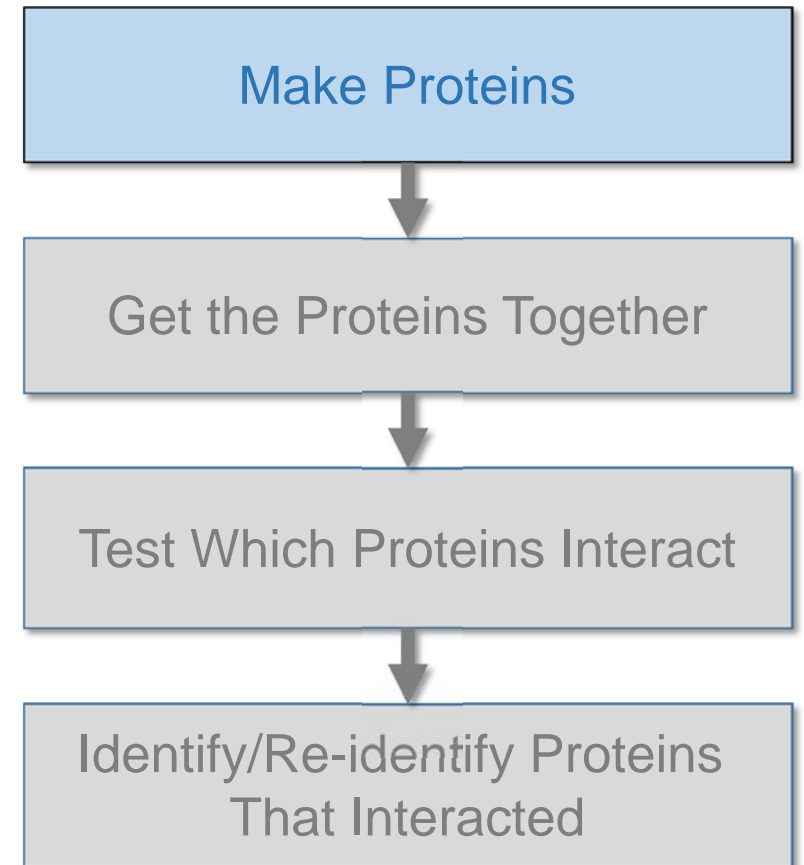
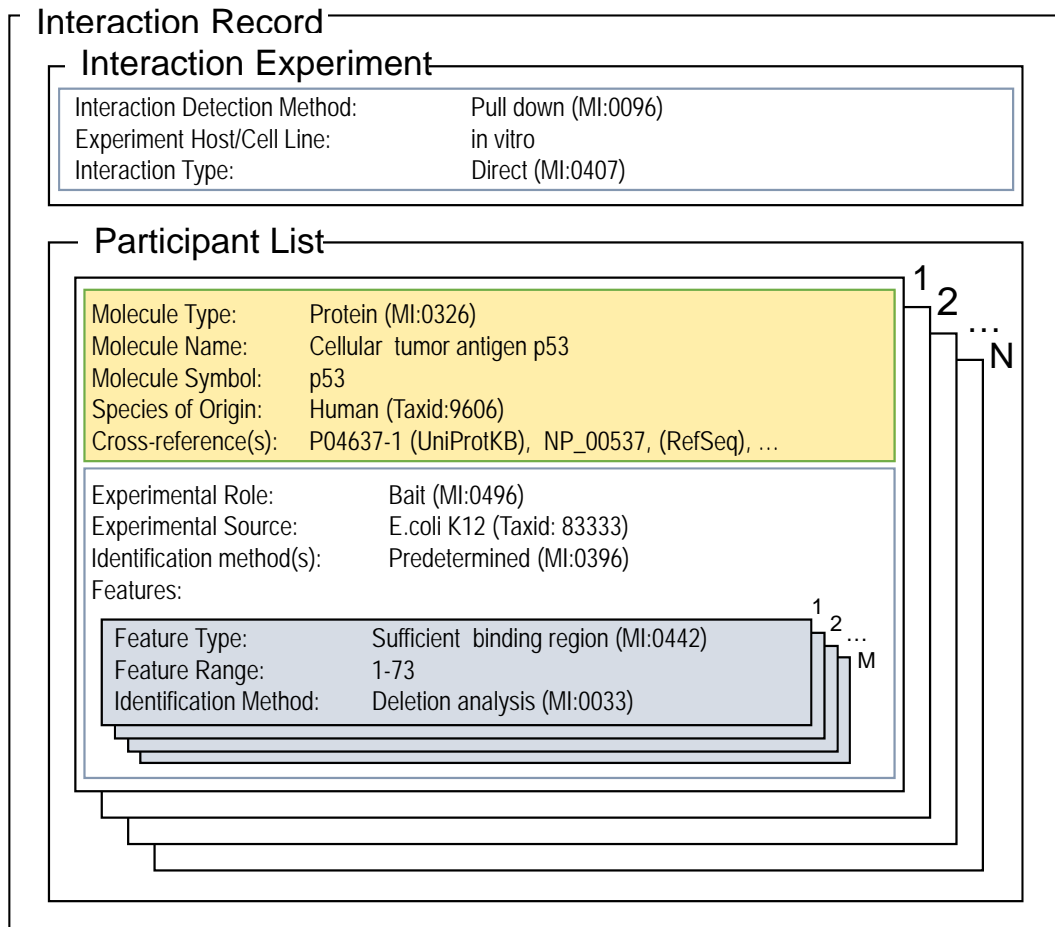
Lukasz Salwinski  
[lukasz@mbi.ucla.edu](mailto:lukasz@mbi.ucla.edu)  
Boyer Hall 205

# Interaction Experiment Flow



- **Make proteins**
  - Where: native vs heterologous host vs in vitro translation vs chemical synthesis
  - How much: native level vs overexpressed
  - Modifications: isoforms, fragments, mutations, PTMs present/absent
- **Get them together**
  - Where: native organism/cell type/tissue/compartiment vs something else
  - When: cell cycle phase/cell state
- **Test which ones interact**
  - Diverse methods can be used to determine that that proteins interact
  - Information that can be inferred from each experiment depends on the method and experimental setup
- **Identify proteins that interact**
  - Identity of some proteins might be known a priori (eg purified, cloned/ tagged bait, etc)
  - Identity and/or state of some proteins might be ambiguous (eg unknown splice form, PTMs)
  - Some molecules participating in the interaction might remain unidentified

# Interaction Experiment Record





# Interaction Database Records

## *Formats*

# Interaction Record Formats

## PSI-MI XML (MIF) format

### Interaction Record

#### Interaction Experiment

|                               |                     |
|-------------------------------|---------------------|
| Interaction Detection Method: | Pull down (MI:0096) |
| Experiment Host/Cell Line:    | in vitro            |
| Interaction Type:             | Direct (MI:0407)    |

#### Participant List

|                     |   |
|---------------------|---|
| Molecule Type:      | Protein (MI:0326)                             |
| Molecule Name:      | Cellular tumor antigen p53                    |
| Molecule Symbol:    | p53   |
| Species of Origin:  | Human (Taxid:9606)                            |
| Cross-reference(s): | P04637-1 (UniProtKB), NP_00537, (RefSeq), ... |

|                           |                           |
|---------------------------|---------------------------|
| Experimental Role:        | Bait (MI:0496)            |
| Experimental Source:      | E.coli K12 (Taxid: 83333) |
| Identification method(s): | Predetermined (MI:0396)   |
| Features:                 |                           |


|                        |                                     |
|------------------------|-------------------------------------|
| Feature Type:          | Sufficient binding region (MI:0442) |
| Feature Range:         | 1-73                                |
| Identification Method: | Deletion analysis (MI:0033)         |



```

+<experimentList></experimentList>
+<interactorList></interactorList>
-<interactionList>
  -<interaction imexId="IM-26392-3" id="2877058">
    -<names>
      <shortLabel>bem1-cla4-1</shortLabel>
    </names>
    +<xref></xref>
    -<experimentRef>2877051</experimentRef>
    </experimentRef>
    -<participantList>
      +<participant id="2877059"></participant>
      -<participant id="2877061">
        +<names></names>
        +<xref></xref>
        <interactorRef>2877055</interactorRef>
        +<biologicalRole></biologicalRole>
        -<experimentalRoleList>
          -<experimentalRole>
            -<names>
              <shortLabel>bait</shortLabel>
              <fullName>bait</fullName>
            </names>
            +<xref></xref>
          </experimentalRole>
          +<experimentalPreparationList></experimentalPreparationList>
          +<featureList>
            +<feature id="2877062"></feature>
            </featureList>
          +<hostOrganismList></hostOrganismList>
        </participant>
      </participantList>
    -<interactionType>
      -<names>
        <shortLabel>physical association</shortLabel>
        <fullName>physical association</fullName>
      </names>
      +<xref></xref>
    </interactionType>
    <modelled>false</modelled>
    <intraMolecular>false</intraMolecular>
    <negative>false</negative>
    +<attributeList></attributeList>
  </interaction>
+<interaction imexId="IM-26392-7" id="2877063"></interaction>
+<interaction imexId="IM-26392-25" id="2877071"></interaction>
+<interaction imexId="IM-26392-27" id="2877078"></interaction>
+<interaction imexId="IM-26392-11" id="2877084"></interaction>
  
```

# MIF/MI Definitions



translating  
the code of life

f t in

JOIN HUPO

Home > About HUPO

## ABOUT HUPO


The Human Proteome Organization (HUPO) is an international scientific organization representing and promoting proteomics through international cooperation and collaborations by fostering the development of new technologies, techniques and training.

### HUPO MISSION STATEMENT

To define and promote proteomics through international cooperation and collaborations by fostering the development of new technologies, techniques and training to better understand human disease.

### Objectives

- Foster global collaboration in major proteomics projects by gathering leading international laboratories in life sciences, bioinformatics, mass spectrometry, systems biology, pathology, and medicine;
- Become the point of contact for proteomics research and commercialization activities worldwide;
- Support large-scale proteomics projects that are aimed at:
  - A mechanistic understanding of fundamental biological processes (often using model




translating  
the code of life

f t in

JOIN HUPO

Home > Initiatives > Proteomics Standards Initiative

## PROTEOMICS STANDARDS INITIATIVE



Website:  
<http://www.psidev.info/>

### Overview of Project

The HUPO Proteomics Standards Initiative (PSI) defines community standards for data representation in proteomics to facilitate data comparison, exchange and verification.

#### PSI Governance

Andy Jones, Chair  
Eric Deutsch, Co-chair  
Sandra Orchard, Co-chair

The main organizational unit of the Proteomics Standards Initiative is the work group. Currently, there are the following work groups:

- CompMS
- Early Career Researcher (ECR) Initiative
- Human Antibody Initiative
- Human Proteome Project
- Initiative on MultiOrganism Proteomes
- Pathology Pillar
- Proteomics Standards Initiative**
- Clinical Proteome Tumor Analysis Consortium (CPTAC)

Q Enter search string

# Interaction Record Formats

PSI-MI XML (MIF) format

## Good

- Stable
  - MIF 2.5 (2007)
  - MIF 3.0 (2018; mostly backward-compatible)
- Database-neutral
  - Developed by HUPO-PSI
- Widely used by data providers
  - IMEx Consortium (DIP, IntAct, MINT,...) – native
  - BioGRID – export
- Expressive enough to describe most of the interaction experiments
  - Multi-protein interactions
  - Protein features (PTMs, mutations, ...)
  - Multiple experimental methods/protein
- Not limited to proteins
  - Nucleic Acids
  - Small Molecules

## Bad

- Overly verbose
  - Redundant open/close tags
  - Several levels of nested elements
  - But compresses quite well – 20x is not that rare
- Does not fit ‘Excel spreadsheet’ paradigm
  - Not too surprising – interaction data is NOT tabular
- Limited set of good quality user-side tools
  - Java JAMI is very versatile but complicated
- Limited support for reporting experiment ambiguities
  - MIF 3.0 provides some support but curation lags
- No support for oligo-/poly-saccharides
  - Limited by the current state of nomenclature
  - No active curation (to my knowledge)
- XML is considered to be hard to work with



# Interaction Record Formats

## PSI-MI tab-delimited (MITAB) format

|    | A                   | B                  | C                       | D                       | E                          | F                                      | G   | H               | I                         | J                                | K                                | L      | M |
|----|---------------------|--------------------|-------------------------|-------------------------|----------------------------|--|---|-----------------|---------------------------|----------------------------------|----------------------------------|--------|---|
| 1  | #ID(s) Interactor A | ID(s) Interactor B | Alt. ID(s) Interactor A | Alt. ID(s) Interactor B | Alias(es) Interactor A     | Alias(es) Interactor B                 | Interaction detection method(s)           | Publicatio      | Publication Identifier(s) | Taxid Interactor A               | Taxid Interactor B               | Source |   |
| 2  | uniprotkb:Q99728-1  | uniprotkb:Q13526   | intact:EBI-21498323     | intact:EBI-714158       | unip psi-mi:q99728-1       | unip psi-mi:pin1_human                 | (display psi-mi:Mi:0007                   | Daza-Marl       | pubmed:31270457           | imex:itaxid:9606(human)          | taxid:9606(psi-mi: "M psi-mi: "M |        |   |
| 3  | uniprotkb:Q99728-1  | uniprotkb:P38398   | intact:EBI-21498323     | intact:EBI-344905       | unip psi-mi:q99728-1       | (display psi-mi:brca1_human            | (display psi-mi: "Mi:0007" (anti tag coim | Daza-Marl       | pubmed:31270457           | imex:itaxid:9606(human)          | taxid:9606(psi-mi: "M psi-mi: "M |        |   |
| 4  | uniprotkb:P38398-1  | uniprotkb:Q13526   | intact:EBI-21498346     | intact:EBI-714158       | unip psi-mi:p38398-1       | (display psi-mi:pin1_human             | (display psi-mi: "Mi:0096" (pull down)    | Daza-Marl       | pubmed:31270457           | imex:itaxid:9606(human)          | taxid:9606(psi-mi: "M psi-mi: "M |        |   |
| 5  | uniprotkb:P38398-1  | uniprotkb:Q06609-1 | intact:EBI-21498346     | intact:EBI-5557721      | psi-mi:p38398-1            | (display psi-mi:q06609-1               | (display psi-mi: "Mi:0096" (pull down)    | Daza-Marl       | pubmed:31270457           | imex:itaxid:9606(human)          | taxid:9606(psi-mi: "M psi-mi: "M |        |   |
| 6  | uniprotkb:P38398-1  | uniprotkb:Q99728   | intact:EBI-21498346     | intact:EBI-473181       | unip psi-mi:p38398-1       | (display psi-mi:barcl1_human           | (display psi-mi: "Mi:0096" (pull down)    | Daza-Marl       | pubmed:31270457           | imex:itaxid:9606(human)          | taxid:9606(psi-mi: "M psi-mi: "M |        |   |
| 7  | uniprotkb:P38398-1  | uniprotkb:Q13526   | intact:EBI-21498346     | intact:EBI-714158       | unip psi-mi:p38398-1       | (display psi-mi:pin1_human             | (display psi-mi: "Mi:0096" (pull down)    | Daza-Marl       | pubmed:31270457           | imex:itaxid:9606(human)          | taxid:9606(psi-mi: "M psi-mi: "M |        |   |
| 8  | uniprotkb:P38398    | uniprotkb:Q13526   | intact:EBI-344905       | unip psi-mi:brca1_human | (display psi-mi:pin1_human | (display psi-mi: "Mi:0096" (pull down) | Daza-Marl                                 | pubmed:31270457 | imex:itaxid:9606(human)   | taxid:9606(psi-mi: "M psi-mi: "M |                                  |        |   |
| 9  | uniprotkb:Q99728-1  | uniprotkb:Q13526   | intact:EBI-21498323     | intact:EBI-714158       | unip psi-mi:q99728-1       | (display psi-mi:pin1_human             | (display psi-mi: "Mi:0096" (pull down)    | Daza-Marl       | pubmed:31270457           | imex:itaxid:9606(human)          | taxid:9606(psi-mi: "M psi-mi: "M |        |   |
| 10 | uniprotkb:Q99728-1  | uniprotkb:Q13526   | intact:EBI-21498323     | intact:EBI-714158       | unip psi-mi:q99728-1       | (display psi-mi:pin1_human             | (display psi-mi: "Mi:0096" (pull down)    | Daza-Marl       | pubmed:31270457           | imex:itaxid:9606(human)          | taxid:9606(psi-mi: "M psi-mi: "M |        |   |
| 11 | uniprotkb:Q13526    | uniprotkb:Q99728-1 | intact:EBI-714158       | unip psi-mi:pin1_human  | (display psi-mi:q99728-1   | (display psi-mi: "Mi:0096" (pull down) | Daza-Marl                                 | pubmed:31270457 | imex:itaxid:9606(human)   | taxid:9606(psi-mi: "M psi-mi: "M |                                  |        |   |
| 12 | uniprotkb:P38398-1  | uniprotkb:Q13526   | intact:EBI-21498346     | intact:EBI-714158       | unip psi-mi:p38398-1       | (display psi-mi:pin1_human             | (display psi-mi: "Mi:0096" (pull down)    | Daza-Marl       | pubmed:31270457           | imex:itaxid:9606(human)          | taxid:9606(psi-mi: "M psi-mi: "M |        |   |
| 13 | uniprotkb:P38398-1  | uniprotkb:Q99728-1 | intact:EBI-21498346     | intact:EBI-21498323     | psi-mi:p38398-1            | (display psi-mi:q99728-1               | (display psi-mi: "Mi:0096" (pull down)    | Daza-Marl       | pubmed:31270457           | imex:itaxid:9606(human)          | taxid:9606(psi-mi: "M psi-mi: "M |        |   |

ID(s) interactor A  
ID(s) interactor B  
Alt. ID(s) interactor A  
Alt. ID(s) interactor B  
Alias(es) interactor A  
Alias(es) interactor B  
Interaction detection method(s)  
Publication 1st author(s)  
Publication Identifier(s)  
Taxid interactor A  
Taxid interactor B  
Interaction type(s)  
Source database(s)  
Interaction identifier(s)  
Confidence value(s)

|                                   |             |
|-----------------------------------|-------------|
| Expansion method(s)               | MIT         |
| Biological role(s) interactor A   |             |
| Biological role(s) interactor B   |             |
| Experimental role(s) interactor A |             |
| Experimental role(s) interactor B |             |
| Type(s) interactor A              |             |
| Type(s) interactor B              |             |
| Xref(s) interactor A              | Annotation  |
| Xref(s) interactor B              | Annotation  |
| Interaction Xref(s)               | Interaction |
|                                   | Host        |

| MITAB 2.6 |                      |
|-----------|----------------------|
| A         | "Mi:0096"(pull down) |
| B         | "Mi:0096"(pull down) |
| Factor A  | "Mi:0096"(pull down) |
| Factor B  | "Mi:0007"(anti tag)  |

|                            |  |
|----------------------------|--|
| Annotation(s) interactor A |  |
| Annotation(s) interactor B |  |
| Interaction annotation(s)  |  |
| Host organism(s)           |  |
| Interaction parameter(s)   |  |
| Creation date              |  |
| Update date                |  |
| Checksum(s) interactor A   |  |
| Checksum(s) interactor B   |  |
| Interaction Checksum(s)    |  |
| Negative                   |  |

Feature(s) interactor A **MITAB 2.7**  
 Feature(s) interactor B  
 Stoichiometry(s) interactor A  
 Stoichiometry(s) interactor B  
 Identification method participant A  
 Identification method participant B

## Column values

| XREF                  | VALUE                    | DESCRIPTION |
|-----------------------|--------------------------|-------------|
| uniprotkb:P12345      | (very important protein) |             |
| wwpdb:1COL wwpdb:3COL |                          |             |
| psi-mi:"MI:0096"      | (pull-down)              |             |
| psi-mi:"MI:0097"      | ("super" tag)            |             |

# Interaction Record Formats

PSI-MI tab-delimited (MITAB) format

## Good

- Easy to read into a spreadsheet
- Supported by third-party libraries

## Bad

- Applicable only to binary interactions
  - Cannot handle multi-protein complexes
- Many columns can be multi-valued
  - Requires custom parsing routines
- Some information originally available in MIF format cannot be stored as MITAB
  - The format is lossy – is, essentially, impossible, to restore fully-featured MIF record from its MITAB representation
- Less stable than MIF
  - MITAB 2.5, 2.7, 2.8 (MIF 2.5 derivatives)
  - MITAB 3.0 (3.0 derivative)

# Interaction Record Formats

BioGRID tab & complex tab formats

## Good

- Easy to read into a spreadsheet

## Bad

- Separate format for binary and multi-protein complexes
- Supports only protein-protein interactions
- Identifies proteins by gene identifiers
  - It works only for organisms that do not splice
- Provides less information than PSI-MI MIF
  - This is because BioGRID extracts less information about interactions that IMEx Consortium databases
- Uses simplified, non-standard CV terms
- See also MITAB deficiencies

# Interaction Record Formats

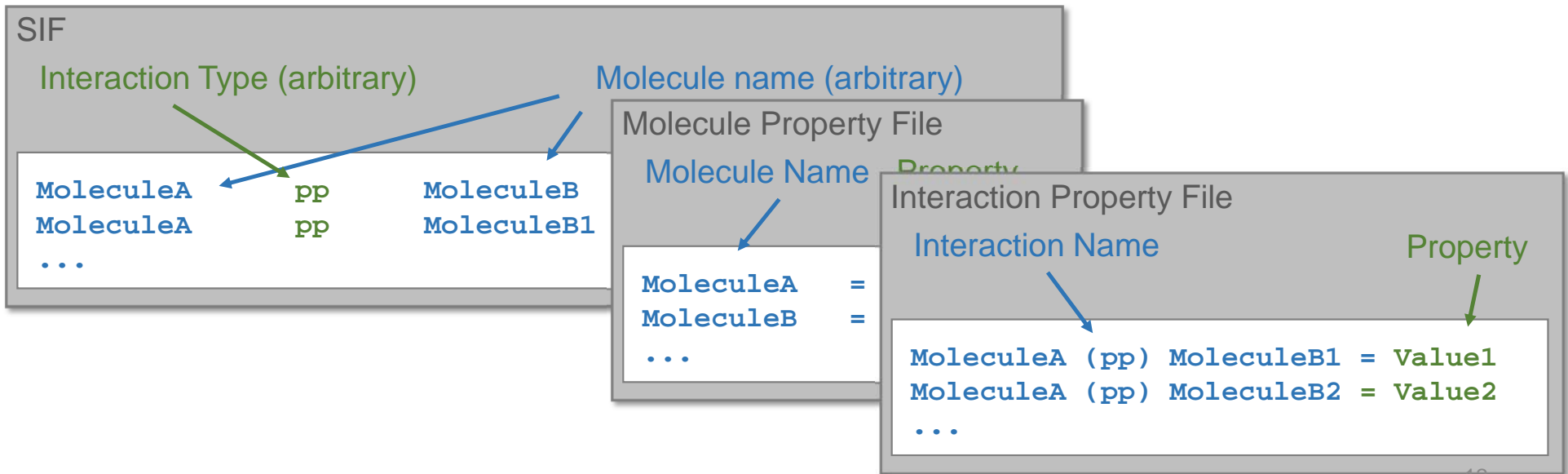
Cytoscape SIF format

**Good**

- Native Cytoscape format
- Simple/easy to prepare (spreadsheet will work)
- Not limited to biological data

**Bad**

- Only binary interactions
- Must be combined with information from additional files in order to provide more detailed information

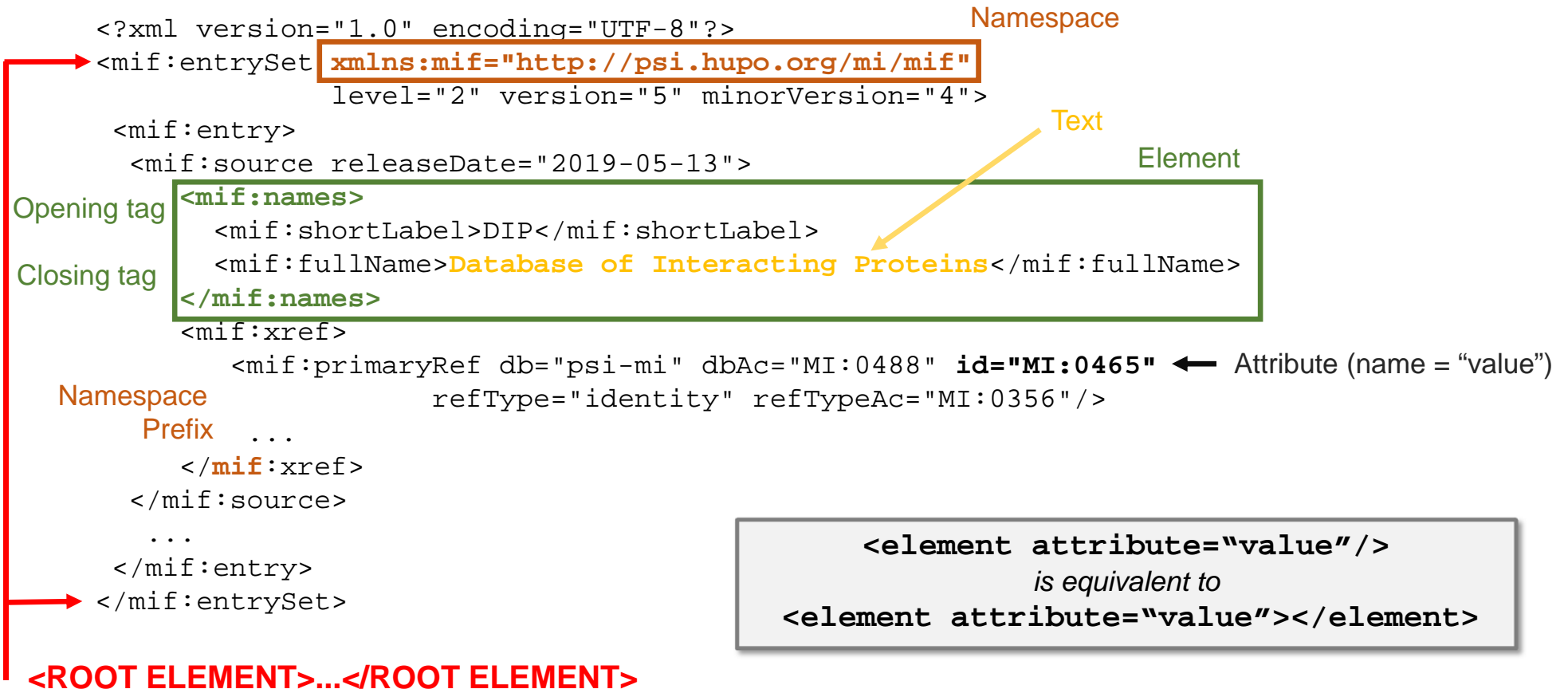


# Interaction Data Processing

## *Topics*

- XML parsing
  - Lxml library
  - Xpath
  - Parsing UniprotKB XML files
- MIF file structure
  - 'compact vs 'expanded' variants
- Binary expansion of multi-protein interactions
  - 'spoke' vs 'matrix' expansion
- Extracting data from MIF files
  - Access relevant/useful information
  - Prepare files ready to import into Cytoscape

# XML File Anatomy



# XML File Parsing

*Lxml library (<https://lxml.de/>)*

```
from lxml import etree
from io import StringIO
```

Only if needed...

```
xml = '<protein acc="P60010"><seq>MKYDDEW...</seq></protein>'
```

```
strDom = etree.fromstring( xml )
```

Parse String

*or*

```
strDom = etree.fromstring( StringIO(xml) )
```

```
fileDom = etree.parse("doc/test.xml")
```

Parse File/URL

*or*

```
fileDom = etree.parse( open("doc/test.xml") )
```

.gz files OK

```
print( etree.tostring( strDom ).decode() )
```

... and back to String

See lxml web site for more options

# XML File Parsing

*lxml library (<https://lxml.de/>)*

```
from lxml import etree
```

```
xml = '<protein acc="P60010"><seq format="fasta">MKYDDEW...</seq></protein>'
```

```
xmlDom = etree.fromstring( xml )
```

```
for child in xmlDom:
```

```
    print( child.tag.decode() )
```

Get element tag

```
    print( child.get("format").decode() )
```

Get attribute

```
    print( child.text.decode() )
```

Get text

```
    print( etree.tostring(child).decode() )
```

Get element as XML

See lxml web site for more options



# XML File Parsing

*lxml Xpath support (<https://lxml.de/xpathxslt.html>)*

```
from lxml import etree
```

```
xml = '<protein acc="P60010"><seq>MKYDDEW...</seq></protein>'
```

```
xmlDom = etree.fromstring( xml )
```

```
root = xmlDom.xpath('/protein')           Get top-level 'protein' element
```

**NOTE: Returns a list !!!**

```
for child in root[0]:
```

```
    print( child.tag.decode() )           Get element tag
```

```
    print( child.get("format").decode() )   Get attribute
```

```
    print( etree.tostring(child).decode() )   Get element as XML
```

See lxml web site for more options

# XML File Parsing

*lxml Xpath support (<https://lxml.de/xpathxslt.html>)*

```
from lxml import etree
```

```
xml = '<protein acc="P60010"><seq>MKYDDEW...</seq></protein>'
```

```
xmlDom = etree.fromstring( xml )
```

```
t1 = xmlDom.xpath('/protein/seq/text()')
```

Get the text of 'seq' elements that are children of the top-level 'protein' element

```
t2 = xmlDom.xpath('//seq/text()')
```

Get the text of ANY 'seq' element

```
e3 = xmlDom.xpath('//protein[@acc="P60010"]/seq')
```

**LIST!!!**  
Get 'seq' elements that are children of 'protein' element with 'acc' attribute equals to 'P60010'

```
e4 = e3[0].xpath('./text()')
```

Get the text of the current element

See <https://en.wikipedia.org/wiki/XPath> and <https://www.w3.org/TR/xpath-10> for more details

# XML File Parsing

*lxml Xpath namespace support (<https://lxml.de/xpathxslt.html>)*

```
from lxml import etree

xml = '''<mif:protein xmlns:mif="http://psi.hupo.org/mi/mif" acc="P60010">
    <mif:seq>MKYDDEW...</mif:seq>
</mif:protein>'''

xmlDom = etree.fromstring( xml )

e = xmlDom.xpath('/m:protein/m:seq',
                 namespaces={'m': 'http://psi.hupo.org/mi/mif'})
```

```
print( e[0].tag.decode() )
```

Get qualified (i.e. with namespace) tag

```
qname = etree.QName(e[0])
```

Split qualified tag into namespace  
and local name

```
print( qname.localname.decode() )
```

```
print( qname.namespace.decode() )
```

# XML File Parsing

*Example: Parsing UniprotKB records*

<https://www.uniprot.org/uniprot/P60010.xml>

```
1 <?xml version='1.0' encoding='UTF-8'?>
2 <uniprot xmlns="http://uniprot.org/uniprot" xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance" xsi:schemaLocation="http://uniprot.org/uniprot/xsd http://uniprot.org/uniprot/xsd">
3   <entry dataset="Swiss-Prot" created="1986-07-21" modified="2019-12-11" version="177">
4     <accession>P60010</accession>
5     <accession>D6VTJ1</accession>
6     <accession>P02579</accession>
7     <accession>Q9F3X6</accession>
8     <accession>Q9F3X7</accession>
9     <name>ACT_YEAST</name>
10    <protein>
11      <recommendedName>
12        <fullName>Actin</fullName>
13      </recommendedName>
14    </protein>
15    <gene>
16      <name type="primary">ACT1</name>
17      <name type="synonym">ABY1</name>
18      <name type="synonym">END7</name>
19      <name type="ordered locus">YFL039C</name>
20    </gene>
21    <organism>
22      <name type="scientific">Saccharomyces cerevisiae (strain ATCC 204508 / S288c)</name>
23      <name type="common">Baker's yeast</name>
24      <dbReference type="NCBI Taxonomy" id="559292"/>
25      <lineage>
26        <taxon>Eukaryota</taxon>
27        <taxon>Fungi</taxon>
28        <taxon>Dikarya</taxon>
29        <taxon>Ascomycota</taxon>
30        <taxon>Saccharomycotina</taxon>
31        <taxon>Saccharomycetes</taxon>
32        <taxon>Saccharomycetales</taxon>
33        <taxon>Saccharomycetaceae</taxon>
34        <taxon>Saccharomyces</taxon>
35      </lineage>
36    </organism>
37    <reference key="1">
38      <citation type="journal article" date="1980" name="Proc. Natl. Acad. Sci. U.S.A." volume="77" first="2546" last="2550">
39        <title>Structure of a split yeast gene: complete nucleotide sequence of the actin gene in Saccharomyces cerevisiae.</title>
40        <authorList>
41          <person name="Gallwitz D."/>
42          <person name="Sures I."/>
43        </authorList>
44        <dbReference type="PubMed" id="6994099"/>
45        <dbReference type="DOI" id="10.1073/pnas.77.5.2546"/>
46      </citation>
47      <scope>NUCLEOTIDE SEQUENCE [GENOMIC DNA]</scope>
48    </reference>
49    <reference key="2">
50      <citation type="journal article" date="1980" name="Proc. Natl. Acad. Sci. U.S.A." volume="77" first="3912" last="3916">
51        <title>Isolation and sequence of the gene for actin in Saccharomyces cerevisiae.</title>
52        <authorList>
53          <person name="Ng R."/>
54          <person name="Abelson J."/>
```

Lots and lots of stuff !!!

- Primary/secondary (aka past) accessions
- Names
  - ❑ Protein Names (including aliases)
  - ❑ Gene Names (including aliases)
- Taxonomy
  - ❑ Species of origin (possibly strain)
  - ❑ Lineage
- Sequence (possibly isoforms)
  - ❑ Molecular weight
- Cross-references
  - ❑ Articles
  - ❑ Gene Ontology Terms
  - ❑ Domains (InterPro, Pfam Prosite)
  - ❑ Structures (PDB)
  - ❑ Interactions
- And more....

# XML File Parsing

Example: Parsing UniprotKB records

```
1 #!/usr/bin/python3
2
3 import sys
4 from lxml import etree
5 from urllib.request import urlopen
6
7 uniprotUrl = "https://www.uniprot.org/uniprot/%s.xml"
8
9 assert len(sys.argv) == 2
10 accession = sys.argv[1]
11
12 accessionUrl = uniprotUrl.replace("%s",accession)
13
14 xmlDoc = etree.parse(urlopen(accessionUrl))
15 print(etree.tostring(xmlDoc).decode())
16
```

uniprotPrint.py

← Needed for https only

## Program flow

- Specify UniprotKB Identifier
- Read Corresponding File
- Find interesting element(s)
- Print results

```
1 #!/usr/bin/python3
2
3 import sys
4 from lxml import etree
5 from urllib.request import urlopen
6
7 uniprotUrl = "https://www.uniprot.org/uniprot/%s.xml"
8 uprotNs = { 'up': 'http://uniprot.org/uniprot' }
9 assert len(sys.argv) == 2
10 accession = sys.argv[1]
11
12 accessionUrl = uprotUrl.replace("%s",accession)
13 uprotTree = etree.parse(urlopen(accessionUrl))
14
15 # find name
16 uprotNameList = uprotTree.xpath("//up:protein/up:fullName/text()",
17                                 namespaces=uprotNs)
18
19 for name in uprotNameList:
20     print("Name: %s" % (name))
21
22 # find accessions
23 uprotAccList = uprotTree.xpath("//up:accession/text()",
24                                 namespaces=uprotNs)
25 accType = "primary"
26 for acc in uprotAccList:
27     print("Accession: %s (%s)" % (acc, accType))
28     accType = "secondary"
29
30 # find gene names
31 uprotGeneNameList = uprotTree.xpath("//up:gene/up:name",
32                                     namespaces=uprotNs)
33
34 for gname in uprotGeneNameList:
35     gnType = gname.xpath("./@type")
36     gnName = gname.xpath("./text()")
37     print("Gene Name: %s (%s)" % (gnName[0], gnType[0]))
38
39 # find sequence
40 uprotSequenceList = uprotTree.xpath("//up:sequence",
41                                     namespaces=uprotNs)
42
43 for seq in uprotSequenceList:
44     seqMass = seq.xpath("./@mass")
45     seqVersion = seq.xpath("./@version")
46     seqStr = seq.xpath("./text()")
47
48 print("Sequence (ver: %s): %s" % (seqVersion[0], seqStr[0]))
49 print("Mass: %6.1fkD" % (int(seqMass[0])/1000))

```

uniprotRead.py

# Interaction Record Formats

## PSI-MI XML (MIF) format

### Interaction Record

#### Interaction Experiment

Interaction Detection Method: Pull down (MI:0096)  
Experiment Host/Cell Line: in vitro  
Interaction Type: Direct (MI:0407)

#### Participant List

Molecule Type: Protein (MI:0326)  
Molecule Name: Cellular tumor antigen p53  
Molecule Symbol: p53  
Species of Origin: Human (Taxid:9606)  
Cross-reference(s): P04637-1 (UniProtKB), NP\_00537, (RefSeq), ...

Experimental Role: Bait (MI:0496)  
Experimental Source: E.coli K12 (Taxid: 83333)  
Identification method(s): Predetermined (MI:0396)  
Features:

Feature Type: Sufficient binding region (MI:0442)  
Feature Range: 1-73  
Identification Method: Deletion analysis (MI:0033)



```
+<experimentList></experimentList>
+<interactorList></interactorList>
-<interactionList>
  -<interaction imexId="IM-26392-3" id="2877058">
    -<names>
      <shortLabel>bem1-cla4-1</shortLabel>
    </names>
    +<xref></xref>
    -<experimentRef>2877051</experimentRef>
    </experimentRef>
    -<participantList>
      +<participant id="2877059"></participant>
      -<participant id="2877061">
        +<names></names>
        +<xref></xref>
        <interactorRef>2877055</interactorRef>
        +<biologicalRole></biologicalRole>
        -<experimentalRoleList>
          -<experimentalRole>
            -<names>
              <shortLabel>bait</shortLabel>
              <fullName>bait</fullName>
            </names>
            +<xref></xref>
            </experimentalRole>
          </experimentalRoleList>
          +<experimentalPreparationList></experimentalPreparationList>
          <featureList>
            +<feature id="2877062"></feature>
            </featureList>
          +<hostOrganismList></hostOrganismList>
        </participant>
      </participantList>
    </interactionType>
    -<names>
      <shortLabel>physical association</shortLabel>
      <fullName>physical association</fullName>
    </names>
    +<xref></xref>
    </interactionType>
    <modelled>false</modelled>
    <intraMolecular>false</intraMolecular>
    <negative>false</negative>
    +<attributeList></attributeList>
  </interaction>
+<interaction imexId="IM-26392-7" id="2877063"></interaction>
+<interaction imexId="IM-26392-25" id="2877071"></interaction>
+<interaction imexId="IM-26392-27" id="2877078"></interaction>
+<interaction imexId="IM-26392-11" id="2877084"></interaction>
```

# Interaction Data Processing

## PSI-MI XML (MIF) file anatomy

Compact MIFs use references to interactors and/or experiments

```
<entry Set>
  <entry>
    <source releaseDate='2019-05-22'>...</source>
    <interactorList>
      <interactor id='1'>...</interactor>
      <interactor id='2'>...</interactor>
      ...
    </interactorList>
    <experimentList>
      <experiment id='1'>...</experiment>
      <experiment id='2'>...</experiment>
      ...
    </experimentList>
    <interactionList>
      <interaction id='1' imexId='IM-123456-1'>
        <experimentList>
          <experimentRef id='1' />
        </experimentList>
        <participantList>
          <participant id='1'>
            <interactorRef id='2' />
            <biologicalRole>...</biologicalRole>
            <experimentalRole>...</experimentalRole>
            ....
          </participant>
          <participant id='2'>...</participant>
          ....
        </participantList>
      </interaction>
      <interaction id='2' imexId='IM-123456-12'>...</interaction>
      ...
    </interactionList>
  </entry>
  ...
</entrySet>
```

Record source (database)

List of interactors

List of experimental methods/protocols

List of interactions

Reference to experimental method/protocol

Reference to interactor

Participant (e.g. protein)

- Describes the state of the molecule as used in the experiment
- Refers to 'interactor' – the reference description of the molecule (e.g. UniprotKB)

# Interaction Data Processing

## PSI-MI XML (MIF) file anatomy

```
<entry Set>
  <entry>
    <source releaseDate='2019-05-22'>...</source>
    <interactionList>
      <interaction id='1' imexId='IM-123456-1'>
        <experimentList>
          <experiment id='1'>...</experiment>
        </experimentList>
        <participantList>
          <participant id='1'>
            <interactor id='2'>...</interactor>
            <biologicalRole>...</biologicalRole>
            <experimentalRole>...</experimentalRole>
            ...
          </participant>
          <participant id='2'>...</participant>
          ...
        </participantList>
      </interaction>
      <interaction id='2' imexId='IM-123456-12'>...</interaction>
      ...
    </interactionList>
  </entry>
  ...
</entrySet>
```

Expanded MIFs describe interactors and experiments within each interaction

Record source (database)

List of interactions

Experimental method/protocol description

Interactor description

Conversion between  
Compact & Expanded  
MIF is lossless !!!

Participant (e.g. protein)

- Describes the state of the molecule as used in the experiment
- Refers to 'interactor' – the reference description of the molecule (e.g. UniprotKB)

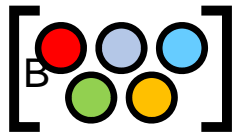


# Interaction Data Processing

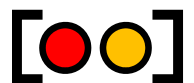
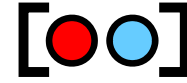
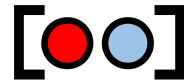
*Binary Expansion*

Binary Expansion is  
is lossy !!!

Association



Physical Associations



Physical Association



Physical Associations



Direct



Direct



# Interaction Data Processing

## PSI-MI XML (MIF) file anatomy

Compact MIFs use references to interactors and/or experiments

```
+<experimentList><experimentList>
+<interactorList><interactorList>
-<interactorList>
+<interaction id="8" imexId="IM-26392-1"><interaction>
-<interaction id="13" imexId="IM-26392-3">
  -<names>
    <shortLabel>bem1-cla4-1</shortLabel>
    <names>
      +<xref><xref>
        -<experimentRef><experimentRef>
        <experimentList>
      -<participantList>
        +<participant id="14"><participant>
          -<participant id="16">
            <interactorRef><interactorRef>
            -<biologicalRole>
              -<names>
                <shortLabel>unspecified role</shortLabel>
                <fullName>unspecified role</fullName>
                <names>
                  +<xref><xref>
                    <biologicalRole>
                    -<experimentalRoleList>
                      -<experimentalRole>
                        -<names>
                          <shortLabel>prey</shortLabel>
                          <fullName>prey</fullName>
                          <names>
                            +<xref><xref>
                              <experimentalRole>
                              <experimentalRoleList>
                                +<experimentalPreparationList><experimentalPreparationList>
                                  -<featureList>
                                    +<feature id="17"><feature>
                                    <featureList>
                                  -<hostOrganismList>
                                    -<hostOrganism ncbiTaxId="559292">
                                      -<names>
                                        <shortLabel>yeast</shortLabel>
                                        <fullName>Saccharomyces cerevisiae</fullName>
                                        <alias type="synonym" typeAc="MI:1041">Baker's yeast</alias>
                                        <names>
                                          <hostOrganism>
                                          <hostOrganismList>
                                        -<attributeList>
                                          <attribute name="comment" nameAc="MI:0612">stoichiometry: 0</attribute>
                                          <attributeList>
                                        </participant>
                                      </participantList>
                                    <interactionType>
                                      -<names>
                                        <shortLabel>physical association</shortLabel>
                                        <fullName>physical association</fullName>
                                        <names>
                                          +<xref><xref>
                                        <interactionType>
                                        <interactionTypeList>
                                          +<attributeList><attributeList>
                                        </interaction>
                                      </interaction>
                                    +<interaction id="18" imexId="IM-26392-2"><interaction>
                                    +<interaction id="23" imexId="IM-26392-21"><interaction>
```

Description of the experimental method (explicit or a reference)

Description of the interactor (explicit or a reference)

Participant (e.g. protein)

- Describes the state of the molecule as used in the experiment
- Refers to 'interactor' – the reference description of the molecule (e.g. UniprotKB)

Features (e.g. mutations)

- Modifications of the molecule relative to the reference state

Interaction

# Interaction Data Processing

## Data structures

### Protein (reference state/interactor)

- ❑ Nested dictionary

### Protein collection

- ❑ Dictionary, using unique id as key

### Interaction evidence

- ❑ Nested dictionary
- ❑ Participants as list
  - ❑ Each refers to interactor through its unique id

### Interaction evidence collection

- ❑ List
  - or
- ❑ Dictionary, using IMEx id as key

```
1
2 prot1 = { "accession" : "P60010",
3           "taxonId" : 4932,
4           "seqStr" : "MDSEVAALVIDNGSGMCKAGFAGDD...",
5           ... }
6
7
8
9
10 prot2 = { "accession": "P60010",
11           "taxon" : { "taxId": 4932,
12                     "sciName" : "S. cerevisiae",
13                     "comName" : "budding yeast" },
14           "seqStr": ['M','D','S',...],
15           "GO": [{ "id" : "GO:0016573", value: "P:histone acetylation"},
16                 { "id" : "GO:0006887", value: "P:exocytosis" }],
17           ... }
18
19
20 protList = [ {}, {}, {} ]
21
22 or
23
24 protDict = { "P60010": prot2, ... }
25
26 Most often easier to work with
27
28
29 evid1 = { "imexId" : "IM-1234-1",
30           "ppantList" : [{ "protId": "P60010",
31                           "expHost": 83333,
32                           "features": [feature1, feature2,...],
33                           "idMethods": ["MI:0123",...],
34                           ...},
35                           { "protId": "P61234",
36                           "expHost": 4932,
37                           "features": [feature3, feature4, ...],
38                           ...},
39                           ...],
40           "intType": { "id": "MI:0912", "value": "direct"},
41           "detectionMth": { "id": "MI:0324", "value": "pull-down"},
42           ... }
43
44 evidList = [ evid1, evid2, evid3, ... ]
45
```

# MIF File Parsing

## Preliminaries/Record Source Information

```
1  #!/usr/bin/python3
2
3  import sys
4  from lxml import etree
5
6  mifNs = { 'm': 'http://psi.hupo.org/mi/mif' }
7  assert len(sys.argv) == 2
8  mifFile = sys.argv[1]
9
10 mifTree = etree.parse( mifFile )
11
12 # find record source
13 mifSrcName = mifTree.xpath( "//m:source/m:shortLabel/text()",
14                             namespaces=mifNs)
15 print( "Record Source: %s" % ( mifSrcName[0] ) )
16
17 # find record release date
18 mifDate = mifTree.xpath( "//m:source/@releaseDate",
19                          namespaces=mifNs)
20 print( "Release Date: %s" % ( mifDate[0] ) )
21
```

mifReadFile.py

## Program flow

- Specify input file location
- Read and parse the input file
- Find interesting element(s)
- Print out/save results

```
1  <?xml version="1.0" encoding="UTF-8"?>
2  <entrySet xmlns="http://psi.hupo.org/mi/mif" xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance">
3    <entry>
4      <source releaseDate="2019-05-13">
5        <names>
6          <shortLabel>DIP</shortLabel>
7          <fullName>Database of Interacting Proteins</fullName>
8        </names>
9        <xref>
10         <primaryRef db="psi-mi" dbAc="MI:0488" id="MI:0465" refType="database">Database of Interacting Proteins</primaryRef>
11         <secondaryRef db="pubmed" dbAc="MI:0446" id="14681454" refType="pubmed">14681454</secondaryRef>
12         <secondaryRef db="intact" dbAc="MI:0469" id="EBI-1579232" refType="intact">EBI-1579232</secondaryRef>
13       </xref>
14       <attributeList>
15         <attribute name="postaladdress">611 Young Drive East; Los Angeles, CA 90095</attribute>
16         <attribute name="url" nameAc="MI:0614">http://dip.doe-mbi.ucla.edu</attribute>
17         <attribute name="contact-email" nameAc="MI:0634">dip@mbi.ucla.edu</attribute>
18       </attributeList>
19     </source>
20     <experimentList>
21       <experimentDescription id="1">DIP</experimentDescription>
22       <experimentDescription id="2">DIP</experimentDescription>
23     </experimentList>
24     <interactorList>
25       <interactor id="3">DIP</interactor>
26       <interactor id="4">DIP</interactor>
27       <interactor id="5">DIP</interactor>
28       <interactor id="6">DIP</interactor>
29       <interactor id="7">DIP</interactor>
30       <interactor id="8">DIP</interactor>
31       <interactor id="9">DIP</interactor>
32       <interactor id="10">DIP</interactor>
33       <interactor id="11">DIP</interactor>
34       <interactor id="12">DIP</interactor>
35       <interactor id="13">DIP</interactor>
36       <interactor id="14">DIP</interactor>
37       <interactor id="15">DIP</interactor>
38       <interactor id="16">DIP</interactor>
39     </interactorList>
40     <interactionList>
41       <interaction id="17">DIP</interaction>
42       <interaction id="33">DIP</interaction>
43     </interactionList>
44   </entry>
45 </entrySet>
```

# MIF File Parsing

## Experimental Method Information

```
22 exByRefId = {}
23
24 #find experiments
25 mifExList = mifTree.xpath( "//m:experimentDescription",
26                             namespaces=mifNs)
27 for ex in mifExList:
28
29     experiment = {}
30
31     # experiment id attribute: unique *within individual file*
32     exId = ex.xpath( "./@id",
33                     namespaces=mifNs)
34     exByRefId[ exId[0] ] = experiment
35
36     #pubmed
37     exPmid = ex.xpath( "./m:bibref/m:xref/*[@db='pubmed']/@id",
38                       namespaces=mifNs)
39     if exPmid:
40         experiment['pmid'] = exPmid[0]
41
42     #experiment host
43     exHost = ex.xpath( "./m:hostOrganismList/m:hostOrganism",
44                       namespaces=mifNs)
45     hostList = []
46     experiment['hostList'] = hostList
47     for host in exHost:
48         curHost = {}
49         hostList.append( curHost )
50
51         hostName = host.xpath( "./m:names/m:shortLabel/text()",
52                               namespaces=mifNs )
53         hostTaxid = host.xpath( "./@ncbiTaxId",
54                                namespaces=mifNs )
55
56         curHost['name'] = hostName[0]
57         curHost['taxid'] = hostTaxid[0]
58
```

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <entrySet xmlns="http://psi.hupo.org/mi/mif" xmlns:xsi="http://www.w3.
3
4 <entry>
5   <source releaseDate="2019-05-13">
6
7   <experimentList>
8     <experimentDescription id="1">
9       <names>
10         <fullName>High-resolution cryo-EM analysis of the yeast ATP
11       </names>
12       <bibref>
13         <xref>
14           <primaryRef db="pubmed" dbAc="MI:0446" id="29650704" refTy
15           <secondaryRef db="intact" dbAc="MI:0469" id="EBI-20595933"
16         </xref>
17       </bibref>
18       <xref>
19         <primaryRef db="pubmed" dbAc="MI:0446" id="29650704" refType
20       </xref>
21     </experimentDescription>
22     <hostOrganismList>
23       <hostOrganism ncbiTaxId="559292">
24         <names>
25           <shortLabel>yeast</shortLabel>
26           <fullName>Saccharomyces cerevisiae</fullName>
27           <alias type="synonym" typeAc="MI:1041">Baker's yeast</al
28         </names>
29       </hostOrganism>
30     </hostOrganismList>
31     <interactionDetectionMethod>
32       <names>
33         <shortLabel>pull down</shortLabel>
34         <fullName>pull down</fullName>
35       </names>
36       <xref>
37         <primaryRef db="psi-mi" dbAc="MI:0488" id="MI:0096" refTyp
38         <secondaryRef db="intact" dbAc="MI:0469" id="EBI-1223" ref
39         <secondaryRef db="pubmed" dbAc="MI:0446" id="14755292" ref
40       </xref>
41     </interactionDetectionMethod>
42     <participantIdentificationMethod>
43       <names>
44         <shortLabel>weight identificat</shortLabel>
45         <fullName>confirmation by molecular weight</fullName>
46         <alias type="synonym" typeAc="MI:1041">weight identificat<
47       </names>
48     </participantIdentificationMethod>
49   </entry>
50 </entrySet>
```

# MIF File Parsing

## Experimental Method Information

```
34 exByRefId[ exId[0] ] = experiment
35
36 #pubmed
37 exPmid = ex.xpath( "./m:bibref/m:xref/*[@db='pubmed']/@id",
38                   namespaces=mifNs)
39 if exPmid:
40     experiment['pmid'] = exPmid[0]
41
42 #experiment host
43 exHost = ex.xpath( "./m:hostOrganismList/m:hostOrganism",
44                   namespaces=mifNs)
45 hostList = []
46 experiment['hostList'] = hostList
47 for host in exHost:
48     curHost = {}
49     hostList.append( curHost )
50
51     hostName = host.xpath( "./m:names/m:shortLabel/text()",
52                           namespaces=mifNs )
53     hostTaxid = host.xpath( "./@ncbiTaxId",
54                            namespaces=mifNs )
55
56     curHost['name'] = hostName[0]
57     curHost['taxid'] = hostTaxid[0]
58
59 # detection method
60 exDetMth = ex.xpath( "./m:interactionDetectionMethod",
61                     namespaces=mifNs)
62 detMthName = exDetMth[0].xpath( "./m:names/m:shortLabel/text()",
63                                 namespaces=mifNs )
64 detMthId = exDetMth[0].xpath( "./m:xref/*[@db='psi-mi']/@id",
65                              namespaces=mifNs )
66
67 experiment['detMeth'] = { 'name': detMthName[0],
68                          'cv': 'psi-mi',
69                          'id': detMthId[0] }
```

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <entrySet xmlns="http://psi.hupo.org/mi/mif" xmlns:xsi="http://www.w3.
3 <entry>
4   <source releaseDate="2019-05-13">
20   <experimentList>
21     <experimentDescription id="1">
22       <names>
23         <fullName>High-resolution cryo-EM analysis of the yeast ATP
24       </names>
25       <bibref>
26         <xref>
27           <primaryRef db="pubmed" dbAc="MI:0446" id="29650704" refTy
28           <secondaryRef db="intact" dbAc="MI:0469" id="EBI-20595933"
29         </xref>
30       </bibref>
31       <xref>
32         <primaryRef db="pubmed" dbAc="MI:0446" id="29650704" refType
33       </xref>
34       <hostOrganismList>
35         <hostOrganism ncbiTaxId="559292">
36           <names>
37             <shortLabel>yeast</shortLabel>
38             <fullName>Saccharomyces cerevisiae</fullName>
39             <alias type="synonym" typeAc="MI:1041">Baker's yeast</al
40           </names>
41         </hostOrganism>
42       </hostOrganismList>
43       <interactionDetectionMethod>
44         <names>
45           <shortLabel>pull down</shortLabel>
46           <fullName>pull down</fullName>
47         </names>
48         <xref>
49           <primaryRef db="psi-mi" dbAc="MI:0488" id="MI:0096" refTy
50           <secondaryRef db="intact" dbAc="MI:0469" id="EBI-1223" re
51           <secondaryRef db="pubmed" dbAc="MI:0446" id="14755292" re
52         </xref>
53       </interactionDetectionMethod>
54       <participantIdentificationMethod>
55         <names>
56           <shortLabel>weight identificat</shortLabel>
57           <fullName>confirmation by molecular weight</fullName>
58           <alias type="synonym" typeAc="MI:1041">weight identificat<
59         </names>
60       </participantIdentificationMethod>
```



# MIF File Parsing

## Interactor (Reference Molecule) Information (I)

```
72 irByRefId = {}
73 irByAccId = {}
74
75 #find interactors
76 mifIrList = mifTree.xpath( "//m:interactor",
77                             namespaces=mifNs)
78 for ir in mifIrList:
79     interactor = {}
80
81     # primary identifier: *shoud* be unique across *all files*
82     primaryIdDB = ir.xpath( "./m:xref/m:primaryRef/@db",
83                             namespaces=mifNs)
84     primaryIdAC = ir.xpath( "./m:xref/m:primaryRef/@id",
85                             namespaces=mifNs)
86
87     interactor['acc'] = primaryIdDB[0]+"."+primaryIdAC[0]
88     irByAccId[ interactor['acc'] ] = interactor
89
90     # interactor id attribute: unique *only* within individual file
91     irId = ir.xpath( "./@id",
92                     namespaces=mifNs)
93     irByRefId[ irId[0] ] = interactor
94
95     # short label
96     label = ir.xpath( "./m:names/m:shortLabel/text()",
97                     namespaces=mifNs)
98     if label:
99         interactor['label'] = label[0]
100     else:
101         interactor['label'] = ''
102
103     # interactor (molecule) type
104     irType = ir.xpath( "./m:interactorType",
105                       namespaces=mifNs )
106     irTypeName = irType[0].xpath( "./m:names/m:shortLabel/text()",
107                                   namespaces=mifNs )
```

```
20 <experimentList>
21 <experimentDescription id="1">
76 <experimentDescription id="2">
133 </experimentList>
134 <interactorList>
135 <interactor id="3">
136 <names>
137 <shortLabel atpg yeast /></shortLabel>
138 <fullName>ATP synthase subunit gamma, mitochondrial</fullName>
139 <alias type="gene name" typeAc="MI:0301">ATP3</alias>
140 <alias type="locus name" typeAc="MI:0305">YBR039W</alias>
141 <alias type="orf name" typeAc="MI:0306">YBR0408</alias>
142 <alias type="gene name synonym" typeAc="MI:0302">F-ATPase ga
143 <alias type="gene name synonym" typeAc="MI:0302">ATP3a</alie
144 <alias type="gene name synonym" typeAc="MI:0302">ATP3b</alie
145 </names>
146 <xref>
147 <primaryRef db="uniprotkb" dbAc="MI:0486" id="P38077" versio
148 <secondaryRef db="uniprotkb" dbAc="MI:0486" id="Q54AF5" vers
149 <secondaryRef db="uniprotkb" dbAc="MI:0486" id="Q76MT6" vers
150 <secondaryRef db="uniprotkb" dbAc="MI:0486" id="D6VQ39" vers
151 <secondaryRef db="intact" dbAc="MI:0469" id="EBI-3271" refTy
152 <secondaryRef db="go" dbAc="MI:0448" id="GO:0005756"/>
153 <secondaryRef db="go" dbAc="MI:0448" id="GO:0046933"/>
154 <secondaryRef db="go" dbAc="MI:0448" id="GO:0015986"/>
155 <secondaryRef db="refseq" dbAc="MI:0481" id="NP_009595.1"/>
156 <secondaryRef db="sgd" dbAc="MI:0484" id="S000000243"/>
157 <secondaryRef db="interpro" dbAc="MI:0449" id="IPR000131"/>
158 <secondaryRef db="interpro" dbAc="MI:0449" id="IPR023632"/>
159 <secondaryRef db="rcsb pdb" dbAc="MI:0460" id="2HLD"/>
160 <secondaryRef db="rcsb pdb" dbAc="MI:0460" id="2WPD"/>
161 <secondaryRef db="rcsb pdb" dbAc="MI:0460" id="2XOK"/>
162 <secondaryRef db="rcsb pdb" dbAc="MI:0460" id="3FKS"/>
163 <secondaryRef db="rcsb pdb" dbAc="MI:0460" id="3OE7"/>
164 <secondaryRef db="rcsb pdb" dbAc="MI:0460" id="3OEE"/>
165 <secondaryRef db="rcsb pdb" dbAc="MI:0460" id="3OEH"/>
166 <secondaryRef db="rcsb pdb" dbAc="MI:0460" id="3OFN"/>
167 <secondaryRef db="rcsb pdb" dbAc="MI:0460" id="3ZRY"/>
168 <secondaryRef db="rcsb pdb" dbAc="MI:0460" id="4B2Q"/>
169 <secondaryRef db="rcsb pdb" dbAc="MI:0460" id="3ZIA"/>
170 <secondaryRef db="dip" dbAc="MI:0465" id="DIP-3035N"/>
171 <secondaryRef db="interpro" dbAc="MI:0449" id="IPR035968"/>
172 <secondaryRef db="rcsb pdb" dbAc="MI:0460" id="6B8H"/>
173 <secondaryRef db="mint" dbAc="MI:0471" id="P38077"/>
174 <secondaryRef db="rcsb pdb" dbAc="MI:0460" id="6CP3"/>
```

# MIF File Parsing

## Interactor (Reference Molecule) Information (II)

```
85 = primaryIdAC = ir.xpath( "./m:xref/m:primaryRef/@id",
86 |                                     namespaces=mifNs)
87
88 interactor['acc'] = primaryIdDB[0] + ":" + primaryIdAC[0]
89 irByAccId[ interactor['acc'] ] = interactor
90
91 # interactor id attribute: unique *only* within individual file
92 = irId = ir.xpath( "./@id",
93 |                 namespaces=mifNs)
94 irByRefId[ irId[0] ] = interactor
95
96 # short label
97 = label = ir.xpath( "./m:names/m:shortLabel/text()",
98 |                 namespaces=mifNs)
99 = if label:
100 |     interactor['label'] = label[0]
101 = else:
102 |     interactor['label'] = ''
103
104 # interactor (molecule) type
105 = irType = ir.xpath( "./m:interactorType",
106 |                 namespaces=mifNs )
107 = irTypeName = irType[0].xpath( "./m:names/m:shortLabel/text()",
108 |                               namespaces=mifNs )
109 = irTypeId = irType[0].xpath( "./m:xref/*[@db='psi-mi']/@id",
110 |                             namespaces=mifNs )
111
112 = interactor['type'] = { 'name': irTypeName[0],
113 |                       'cv': 'psi-mi',
114 |                       'id': irTypeId[0] }
115
116 # gene name
117 = gene = ir.xpath( "./m:names/m:alias[@type='gene name']/text()",
118 |                 namespaces=mifNs)
119 = if gene:
120 |     interactor['gene'] = gene[0]
121
```

```
134 = <interactorList>
135 = <interactor id="3">
136 = <names>
137 | <shortLabel>atpg_yeast</shortLabel>
138 | <fullName>ATP synthase subunit gamma, mitochondrial</fullName>
139 | <alias type="gene name" typeAc="MI:0301">ATP3</alias>
140 | <alias type="locus name" typeAc="MI:0305">YBR039W</alias>
141 | <alias type="orf name" typeAc="MI:0306">YBR0408</alias>
142 | <alias type="gene name synonym" typeAc="MI:0302">F-ATPase ga
143 | <alias type="gene name synonym" typeAc="MI:0302">ATP3a</alie
144 | <alias type="gene name synonym" typeAc="MI:0302">ATP3b</alis
145 = </names>
146 = <xref>
181 = <interactorType>
182 = <names>
183 | <shortLabel>protein</shortLabel>
184 | <fullName>protein</fullName>
185 = </names>
186 = <xref>
187 | <primaryRef db="psi-mi" dbAc="MI:0488" id="MI:0326" refTy
188 | <secondaryRef db="intact" dbAc="MI:0469" id="EBI-619654"
189 | <secondaryRef db="pubmed" dbAc="MI:0446" id="14755292" re
190 | <secondaryRef db="so" dbAc="MI:0601" id="so:0000358" refT
191 = </xref>
192 = </interactorType>
193 = <organism ncbiTaxId="559292">
194 = <names>
195 | <shortLabel>yeast</shortLabel>
196 | <fullName>Saccharomyces cerevisiae</fullName>
197 | <alias type="synonym" typeAc="MI:1041">Baker's yeast</alie
198 = </names>
199 = </organism>
200 = <sequence>MLSRIVSNNATRSVMCHQAQVGILYKTNFVRTYATLKEVEMRLKSIKIEKI
201 = <attributeList>
202 | <attribute name="crc64">ADC71F3C1E0CDF91</attribute>
203 | <attribute name="crc64">ADC71F3C1E0CDF91</attribute>
204 = </attributeList>
205 = </interactor>
206 = <interactor id="4">
207 = <names>
208 | <shortLabel>atp5e_yeast</shortLabel>
209 | <fullName>ATP synthase subunit epsilon, mitochondrial</fullN
210 | <alias type="gene name" typeAc="MI:0301">ATP15</alias>
211 | <alias type="orf name" typeAc="MI:0306">P0345</alias>
212 | <alias type="locus name" typeAc="MI:0305">YPL271W</alias>
```



# MIF File Parsing

## Interaction Evidence Information (I)

```
123 evidList = []
124
125 #find evidence
126 mifEvList = mifTree.xpath( "//m:interactionList/m:interaction",
127                             namespaces=mifNs)
128 for ev in mifEvList:
129
130     evid = {}
131     evidList.append( evid )
132     expt = None
133
134     # find experiment
135     evExpRef = ev.xpath( "../m:experimentRef/text()",
136                         namespaces=mifNs )
137     if evExpRef: # compact MIF
138         expt = exByRefId[ evExpRef[0] ]
139     else: # expanded MIF
140         evExpId = ev.xpath( "../m:experimentDescription/@id",
141                             namespaces=mifNs )
142         if evExpId:
143             expt = exByRefId[ evExpId[0] ]
144
145     evid['expt'] = expt
146
147     # find interaction type
148     evITname = ev.xpath( "../m:interactionType//m:shortLabel/text()",
149                         namespaces=mifNs )
150     evITid = ev.xpath( "../m:interactionType/m:xref/*[@db='psi-mi']/@id",
151                       namespaces=mifNs )
152
153     evid['IntType'] = { 'name': evITname[0],
154                       'cv': 'psi-mi',
155                       'id': evITid[0] }
156
157     ptList = []
158     evid['ptList'] = ptList
159
```

```
1023 <interactionList>
1024   <interaction id="17">
1025     <names>
1026       <shortLabel>atp14-atp18-1</shortLabel>
1027     </names>
1028     <xref>
1029       <primaryRef db="wwpdb" dbAc="MI:0805" id="6CP3" refType="id"
1030       <secondaryRef db="wwpdb" dbAc="MI:0805" id="6CP5" refType="
1031       <secondaryRef db="emdb" dbAc="MI:0936" id="EMD-7548" refTyp
1032       <secondaryRef db="emdb" dbAc="MI:0936" id="EMD-7549" refTyp
1033       <secondaryRef db="emdb" dbAc="MI:0936" id="EMD-7546" refTyp
1034       <secondaryRef db="emdb" dbAc="MI:0936" id="EMD-7547" refTyp
1035       <secondaryRef db="wwpdb" dbAc="MI:0805" id="6CP6" refType="
1036       <secondaryRef db="wwpdb" dbAc="MI:0805" id="6CP7" refType="
1037       <secondaryRef db="intact" dbAc="MI:0469" id="EBI-20595952"
1038     </xref>
1039     <experimentList>
1040       <experimentRef>2</experimentRef>
1041     </experimentList>
1042     <participantList>
1043       <participant id="18">
1044       <participant id="19">
1045       <participant id="20">
1046       <participant id="21">
1047       <participant id="22">
1048       <participant id="23">
1049       <participant id="24">
1050       <participant id="25">
1051       <participant id="26">
1052       <participant id="28">
1053       <participant id="29">
1054       <participant id="30">
1055       <participant id="31">
1056       <participant id="32">
1057     </participantList>
1058     <interactionType>
1059       <names>
1060         <shortLabel>physical association</shortLabel>
1061         <fullName>physical association</fullName>
1062       </names>
1063       <xref>
1064         <primaryRef db="psi-mi" dbAc="MI:0488" id="MI:0915" refTy
1065         <secondaryRef db="intact" dbAc="MI:0469" id="EBI-1813.47" re
1066         <secondaryRef db="pubmed" dbAc="MI:0446" id="14755292" re
1067       </xref>
1068     </interactionType>
1069   </interaction>
1070 </interactionList>
```

# MIF File Parsing

## Interaction Evidence Information (II)

```
166 ptList = []
167 evid['ptList'] = ptList
168
169 # find participants
170 evParts = ev.xpath( "../m:participantList/m:participant",
171                     namespaces=mifNs )
172
173 for part in evParts:
174     cpart = {}
175     ptList.append(cpart)
176
177     # find interactor
178     interactor = None
179     intRefId = part.xpath( "../m:interactorRef/text()",
180                           namespaces=mifNs )
181     if intRefId: # compact MIF
182         interactor = irById[ intRefId[0] ]
183     else: # expanded MIF
184         intId = part.xpath( "../m:interactor/@id",
185                             namespaces=mifNs )
186         if intId:
187             interactor = irById[ intId[0] ]
188     cpart['interactor'] = interactor
189
190 # experimental/biological role
191 # (bait, prey, ancillary, enzyme, substrate, etc.)
192 partRoleList = []
193 cpart['roleList'] = partRoleList
194 roleList = part.xpath( "../*[local-name()='experimentalRole'
195                           ' local-name()='biologicalRole' ]",
196                         namespaces=mifNs )
197 for r in roleList:
198     roleName = r.xpath( "../m:names/m:shortLabel/text()",
199                         namespaces=mifNs )
200     roleId = r.xpath( "../m:xref/*[@db='psi-mi']/@id",
201                      namespaces=mifNs )
202
```

```
1042 <participantList>
1043 <participant id="18">
1044   <interactorRef>15</interactorRef>
1045   <biologicalRole>
1046     <names>
1047       <shortLabel>unspecified role</shortLabel>
1048       <fullName>unspecified role</fullName>
1049     </names>
1050     <xref>
1051       <primaryRef db="psi-mi" dbAc="MI:0488" id="MI:0499" r
1052       <secondaryRef db="intact" dbAc="MI:0469" id="EBI-7778
1053       <secondaryRef db="pubmed" dbAc="MI:0446" id="14755292
1054     </xref>
1055   </biologicalRole>
1056   <experimentalRoleList>
1057     <experimentalRole>
1058       <names>
1059         <shortLabel>neutral component</shortLabel>
1060         <fullName>neutral component</fullName>
1061       </names>
1062       <xref>
1063         <primaryRef db="psi-mi" dbAc="MI:0488" id="MI:0497"
1064         <secondaryRef db="intact" dbAc="MI:0469" id="EBI-55
1065         <secondaryRef db="pubmed" dbAc="MI:0446" id="147552
1066       </xref>
1067     </experimentalRole>
1068   </experimentalRoleList>
1069   <experimentalPreparationList>
1070     <experimentalPreparation>
1071       <names>
1072         <shortLabel>purified</shortLabel>
1073         <fullName>purified</fullName>
1074       </names>
1075       <xref>
1076         <primaryRef db="psi-mi" dbAc="MI:0488" id="MI:0350"
1077         <secondaryRef db="intact" dbAc="MI:0469" id="EBI-15
1078         <secondaryRef db="pubmed" dbAc="MI:0446" id="147552
1079       </xref>
1080     </experimentalPreparation>
1081   </experimentalPreparationList>
1082   <hostOrganismList>
1083     <hostOrganism ncbiTaxId="559292">
1084       <names>
1085         <shortLabel>yeast</shortLabel>
1086         <fullName>Saccharomyces cerevisiae</fullName>

```

# MIF File Parsing

## Interaction Evidence Information (III)

```
164 for part in evParts:
165     cpart = {}
166     ptList.append(cpart)
167
168     # find interactor
169     interactor = None
170     intRefId = part.xpath( "./m:interactorRef/text()",
171                             namespaces=mifNs )
172     if intRefId: # compact MIF
173         interactor = irByRefId[ intRefId[0] ]
174     else: # expanded MIF
175         intId = part.xpath( "./m:interactor/@id",
176                             namespaces=mifNs )
177         if intId:
178             interactor = irByRefId[ intId[0] ]
179         cpart['interactor'] = interactor
180
181     # experimental/biological role
182     # (bait, prey, ancillary, enzyme, substrate, etc.)
183     partRoleList = []
184     cpart['roleList'] = partRoleList
185     roleList = part.xpath( ".*[local-name()='experimentalRole'
186                             local-name()='biologicalRole' ]",
187                             namespaces=mifNs )
188     for r in roleList:
189         roleName = r.xpath( "./m:names/m:shortLabel/text()",
190                             namespaces=mifNs )
191         roleId = r.xpath( "./m:xref/*[@db='psi-mi']/@id",
192                             namespaces=mifNs )
193
194         if roleName[0] != 'unspecified role':
195             role = { 'name': roleName[0],
196                     'cv': 'psi-mi',
197                     'id': roleId[0] }
198             partRoleList.append(role)
199
200 import json
201 print( json.dumps( evidList, sort_keys=True, indent=4) )
```

mifFileRead.py

```
1042 <participantList>
1043 <participant id="18">
1044 <interactorRef>15</interactorRef>
1045 <biologicalRole>
1046 <names>
1047 <shortLabel>unspecified role</shortLabel>
1048 <fullName>unspecified role</fullName>
1049 </names>
1050 <xref>
1051 <primaryRef db="psi-mi" dbAc="MI:0488" id="MI:0499">
1052 <secondaryRef db="intact" dbAc="MI:0469" id="EBI-7778">
1053 <secondaryRef db="pubmed" dbAc="MI:0446" id="14755292">
1054 </xref>
1055 </biologicalRole>
1056 <experimentalRoleList>
1057 <experimentalRole>
1058 <names>
1059 <shortLabel>neutral component</shortLabel>
1060 <fullName>neutral component</fullName>
1061 </names>
1062 <xref>
1063 <primaryRef db="psi-mi" dbAc="MI:0488" id="MI:0497">
1064 <secondaryRef db="intact" dbAc="MI:0469" id="EBI-55">
1065 <secondaryRef db="pubmed" dbAc="MI:0446" id="14755292">
1066 </xref>
1067 </experimentalRole>
1068 </experimentalRoleList>
1069 <experimentalPreparationList>
1070 <experimentalPreparation>
1071 <names>
1072 <shortLabel>purified</shortLabel>
1073 <fullName>purified</fullName>
1074 </names>
1075 <xref>
1076 <primaryRef db="psi-mi" dbAc="MI:0488" id="MI:0350">
1077 <secondaryRef db="intact" dbAc="MI:0469" id="EBI-15">
1078 <secondaryRef db="pubmed" dbAc="MI:0446" id="14755292">
1079 </xref>
1080 </experimentalPreparation>
1081 </experimentalPreparationList>
1082 <hostOrganismList>
1083 <hostOrganism ncbiTaxId="559292">
1084 <names>
1085 <shortLabel>yeast</shortLabel>
1086 <fullName>Saccharomyces cerevisiae</fullName>
```

# MIF File Parsing

## Merging Interaction Data (I)

mifFileMerge.py

```
1  #!/usr/bin/python3
2
3  import sys
4  import json
5
6  from os import walk
7  from os.path import join
8
9  from lxml import etree
10
11  mifNs = { 'm': 'http://psi.hupo.org/mi/mif' }
12
13  def getSource( mifData, mifTree ):
14
15      if 'sourceList' not in mifData:
16          mifData['sourceList'] = []
17
18      # find record source
19      mifSrcName = mifTree.xpath( "//m:source//m:shortLabel/text()",
20                                  namespaces=mifNs )
21      # find record release date
22      mifDate = mifTree.xpath( "//m:source/@releaseDate",
23                              namespaces=mifNs )
24
25      mifData['sourceList'].append( { "database": mifSrcName[0],
26                                     "date": mifDate[0] } )
27      return mifData
28
29
30  def getExperiments( mifData, mifTree ):
31
32      mifData['exByRefId'] = {}
33
34      #find experiments
35      mifExList = mifTree.xpath( "//m:experimentDescription",
36                                 namespaces=mifNs )
37      for ex in mifExList:
```

```
1  <?xml version="1.0" encoding="UTF-8"?>
2  <entrySet xmlns="http://psi.hupo.org/mi/mif" xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance">
3      <entry>
4          <source releaseDate="2019-05-13">
5              <names>
6                  <shortLabel>DIP</shortLabel>
7                  <fullName>Database of Interacting Proteins</fullName>
8              </names>
9              <xref>
10                  <primaryRef db="psi-mi" dbAc="MI:0488" id="MI:0465" refType="id"></primaryRef>
11                  <secondaryRef db="pubmed" dbAc="MI:0446" id="14681454" refType="id"></secondaryRef>
12                  <secondaryRef db="intact" dbAc="MI:0469" id="EBI-1579232" refType="id"></secondaryRef>
13              </xref>
14              <attributeList>
15                  <attribute name="postaladdress">611 Young Drive East; Los Angeles, CA 90095-1606</attribute>
16                  <attribute name="url" nameAc="MI:0614">http://dip.doe-mbi.ucla.edu</attribute>
17                  <attribute name="contact-email" nameAc="MI:0634">dip@mbi.ucla.edu</attribute>
18              </attributeList>
19          </source>
20          <experimentList>
21              <experimentDescription id="1"></experimentDescription>
22              <experimentDescription id="2"></experimentDescription>
23          </experimentList>
24          <interactorList>
25              <interactor id="3"></interactor>
26              <interactor id="4"></interactor>
27              <interactor id="5"></interactor>
28              <interactor id="6"></interactor>
29              <interactor id="7"></interactor>
30              <interactor id="8"></interactor>
31              <interactor id="9"></interactor>
32              <interactor id="10"></interactor>
33              <interactor id="11"></interactor>
34              <interactor id="12"></interactor>
35              <interactor id="13"></interactor>
36              <interactor id="14"></interactor>
37              <interactor id="15"></interactor>
38              <interactor id="16"></interactor>
39          </interactorList>
40          <interactionList>
41              <interaction id="17"></interaction>
42              <interaction id="33"></interaction>
43          </interactionList>
44      </entry>
45  </entrySet>
```



# MIF File Parsing

## Merging Interaction Data (II)

```
428 mifData = {}
429
430 assert len(sys.argv) == 4
431 mydir = sys.argv[1]
432 mymth = sys.argv[2]
433 myout = sys.argv[3]
434
435 myfiles = []
436 for (dirpath, dirnames, filenames) in walk( mydir ):
437     myfiles.extend(filenames)
438     break
439
440 for f in myfiles:
441     mifTree = etree.parse( join( mydir, f ) )
442
443     getSource( mifData, mifTree )
444     getExperiments( mifData, mifTree )
445     getInteractors( mifData, mifTree )
446     getEvidence( mifData, mifTree )
447     dropInternalRefs( mifData )
448
449 #-----
450 # operations
451 #-----
452
453 if myout == 'json':
454     print( toJson( mifData, indent=1 ) )
455
456 if mymth in ['spoke', 'matrix']:
457     network = buildNetwork( mifData, expand=mymth
458
459 if myout == 'sif':
460     print( toSif( mifData, network ) )
461 elif myout == 'nodeprop':
462     print( toNodePoperties( mifData, network ) )
463 elif myout == 'edgeprop':
464     print( toEdgePoperties( mifData, network ) )
```

mifFileMerge.py

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <entrySet xmlns="http://psi.hupo.org/mi/mif" xmlns:xsi="http://www.w3.o
3 <entry>
4   <source releaseDate="2019-05-13">
5     <names>
6       <shortLabel>DIP</shortLabel>
7       <fullName>Database of Interacting Proteins</fullName>
8     </names>
9     <xref>
10       <primaryRef db="psi-mi" dbAc="MI:0488" id="MI:0465" refType="id
11       <secondaryRef db="pubmed" dbAc="MI:0446" id="14681454" refType=
12       <secondaryRef db="intact" dbAc="MI:0469" id="EBI-1579232" refTy
13     </xref>
14     <attributeList>
15       <attribute name="postaladdress">611 Young Drive East; Los Angel
16       <attribute name="url" nameAc="MI:0614">http://dip.doe-mbi.ucla.
17       <attribute name="contact-email" nameAc="MI:0634">dip@mbi.ucla.e
18     </attributeList>
19   </source>
20   <experimentList>
21     <experimentDescription id="1">
22       <experimentDescription id="2">
23     </experimentList>
24   </experimentList>
25   <interactorList>
```

stackoverflow Products Customers Use cases

How do I list all files of a directory

Log in Sign up

Home

PUBLIC

Stack Overflow

Tags

Users

TEAMS

Free 30 Day Trial

What's new?

os.listdir() will get you everything that's in a directory - files and directories.

3868

If you want just files, you could either filter this down using [os.path](#):

```
from os import listdir
from os.path import isfile, join
onlyfiles = [f for f in listdir(mypath) if isfile(join(mypath, f))]
```

or you could use [os.walk\(\)](#) which will yield two lists for each directory it visits - splitting into files and dirs for you. If you only want the top directory you can just break the first time it yields

```
from os import walk

f = []
for (dirpath, dirnames, filenames) in walk(mypath):
    f.extend(filenames)
    break
```

share improve this answer

edited Jun 5 '19 at 3:10

answered Jul 8 '10 at 21:01

cs95 189k • 38 • 294 • 358

pycruft 41.8k • 1 • 14 • 10

Linked

- 16 How to get only files in directory?
- 10 Reading all files in all directories
- 11 Python3 list files from particular directory
- 14 python : get list all \*.txt files in a directory
- 5 print the directory for files in Python
- 4 How to list directory using Python
- 1 How can I automatically open all text files in a given folder?
- 2 How do i use Linux terminal commands like CD and LS?
- 2 How to get file path + file name into a list?
- 2 Automatically creating a list in Python

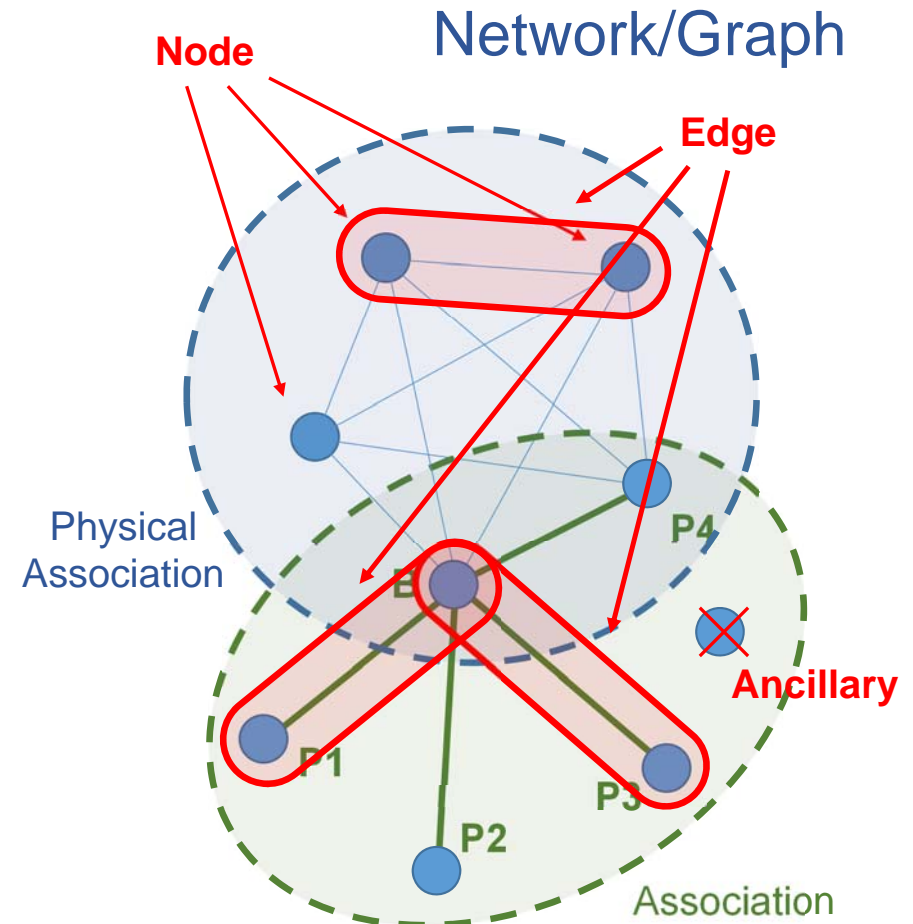
see more linked questions...

## Merging Interaction Data

```

254
255 def buildNetwork( mifData, expand=None):
256
257     # expand:
258     # None - binary edges only
259     # spoke - bait/prey (spoke) expansion
260     # matrix - spoke + matrix of physical and below, ignoring ancillary
261
262     # network data (note: participant details are losts !!!)
263     # key: (acc1,acc2);
264     # value: {'evidList': [ev1,ev2,...], 'intType':[...]}
265     edges = {}
266
267     for ev in mifData['evidList']:
268         bait = None
269         preyDict = {}
270         partDict = {}
271         for p in ev['partList']:
272             skip = False
273             for r in p['roleList']:
274                 if r['name'] == 'bait':
275                     bait = p['interactor']['acc']
276                 if r['name'] == 'prey':
277                     if p['interactor']['acc'] in preyDict:
278                         preyDict[ p['interactor']['acc'] ] += 1
279                     else:
280                         preyDict[ p['interactor']['acc'] ] = 1
281
282                 if r['name'] == 'ancillary':
283                     skip = True
284             if not skip:
285                 if p['interactor']['acc'] in partDict:
286                     partDict[ p['interactor']['acc'] ] += 1
287                 else:
288                     partDict[ p['interactor']['acc'] ] = 1
289
290     if len(partDict) == 2:
291         (acc1,acc2) = sorted(partDict.keys())

```

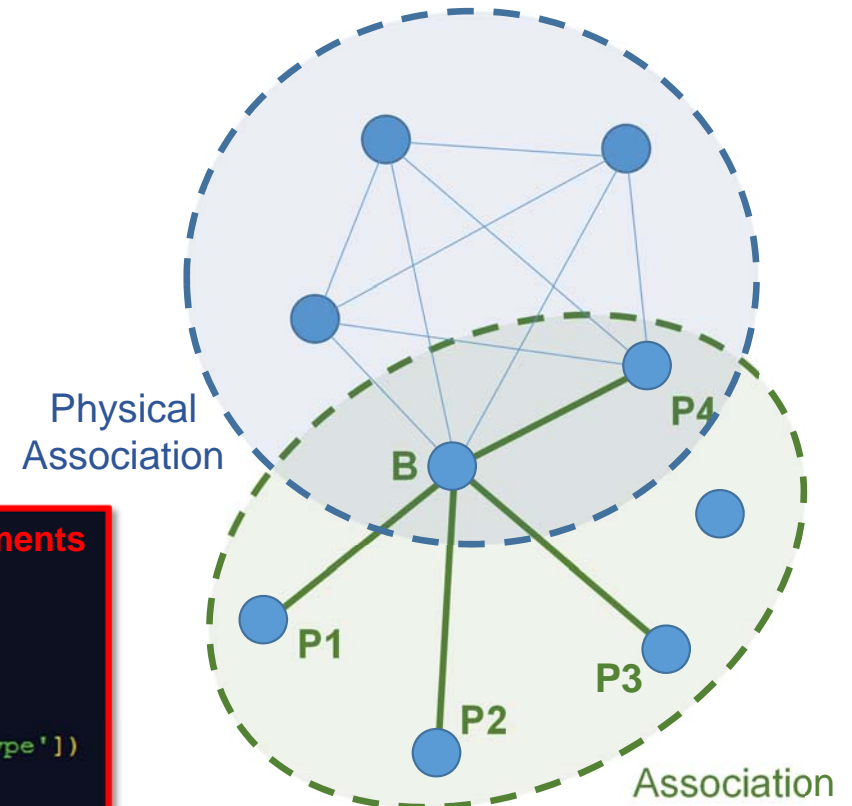


# MIF File Parsing

## Merging Interaction Data

```
254
255 def buildNetwork( mifData, expand=None):
256
257     # expand:
258     # None - binary edges only
259     # spoke - bait/prey (spoke) expansion
260     # matrix - spoke + matrix of physical and below, ignoring ancillar
261
262     # network data (note: participant details are losts !!!)
263     # key: (acc1,acc2);
264     # value: {'evidList': [ev1,ev2,...], 'intType':[...]}
265     edges = {}
266
267     for ev in mifData['evidList']:
268         bait = None
269         preyDict = {}
270         partDict = {}
271         for p in ev['partList']:
272             skip = False
273             for r in p['roleList']:
274
275                 if len(partDict) == 2:
276                     (part1, part2) = ( partDict.keys() )
277                     if part1 > part2:
278                         (part2, part1) = (part1, part2)
279                     if (part1, part2) in edges:
280                         edges[(part1, part2)]['evidList'].append(ev )
281                         edges[(part1, part2)]['intType'].append(ev['intType'])
282                     elif part1 != part2:
283                         edges[(part1, part2)] = {'evidList': [ ev ],
284                                                     'intType': [ev['intType']]}
285
286                 if len(partDict) == 2:
287                     (part1, part2) = ( partDict.keys() )
```

## Network/Graph



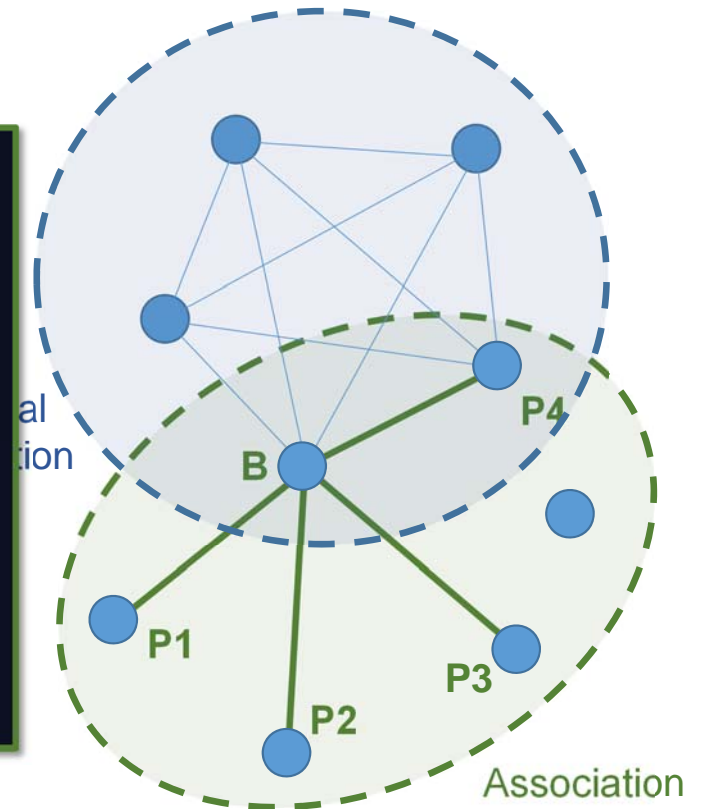


# MIF File Parsing

## Merging Interaction Data

```
254
255 def buildNetwork( mifData, expand=None):
256
257     # expand:
258     # None - binary edges only
259     # spoke - bait/prey (spoke) expansion
260
261     elif len(partDict) > 2:
262         if expand in ['spoke', 'matrix']:
263             if bait is not None:
264                 for prey in preyDict.keys():
265
266
267                 itype = {'cv': 'psi-mi', 'id': 'MI:0915',
268                         'name': 'physical association' }
269
270
271                 if bait > prey:
272                     (part1, part2) = (prey, bait)
273                 else:
274                     (part1, part2) = (bait, prey)
275                 if (part1, part2) in edges:
276                     edges[(part1, part2)]['evidList'].append( ev )
277                     edges[(part1, part2)]['intType'].append( itype )
278                 elif part1 != part2:
279                     edges[ (part1, part2) ] = { 'evidList': [ ev ],
280                                                 'intType': [ itype ] }
281
282
283                 skip = True
284                 if not skip:
285                     if p['interactor']['acc'] in partDict:
286                         partDict[ p['interactor']['acc'] ] += 1
287                     else:
288                         partDict[ p['interactor']['acc'] ] = 1
289
290                 if len(partDict) == 2:
```

## Network/Graph



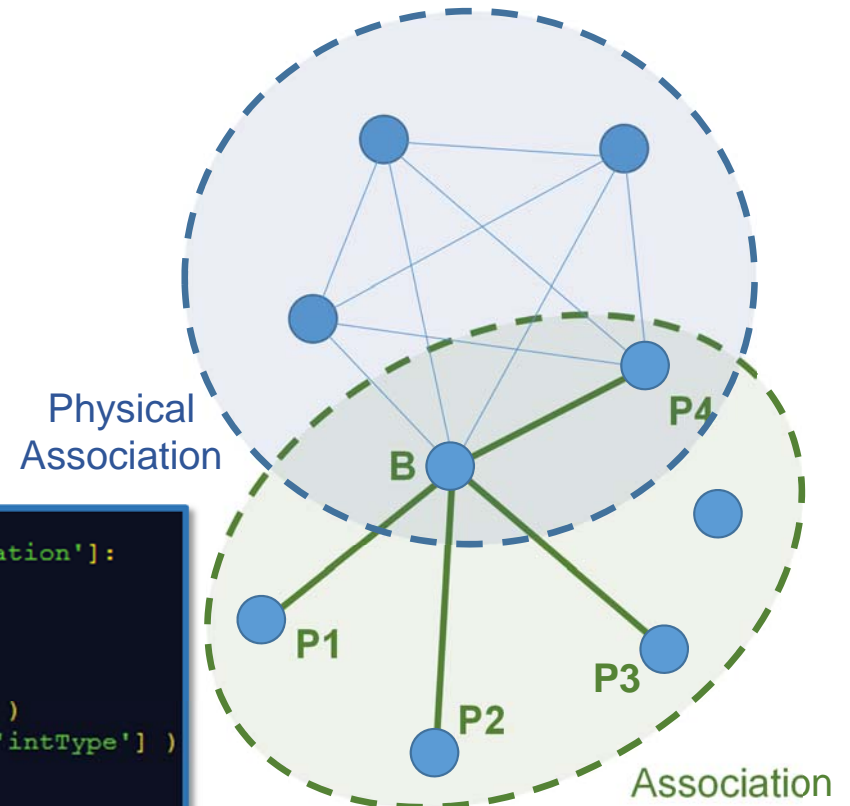


# MIF File Parsing

## Merging Interaction Data

```
254
255 def buildNetwork( mifData, expand=None):
256
257     # expand:
258     # None - binary edges only
259     # spoke - bait/prey (spoke) expansion
260     # matrix - spoke + matrix of physical and below, ignoring ancillar
261
262     # network data (note: participant details are losts !!!)
263     # key: (acc1,acc2);
264     # value: {'evidList': [ev1,ev2,...], 'intType':[...]}
265     edges = {}
266
267     for ev in mifData['evidList']:
268         bait = None
269         preyDict = {}
270         partDict = {}
271         for p in ev['partList']:
272             skip = False
273             for r in p['roleList']:
274
275 if expand in ['matrix'] and ev['intType']['name'] not in ['association']:
276     for part1 in partDict.keys():
277         for part2 in partDict.keys():
278             if part2 > part1:
279                 if (part1, part2) in edges:
280                     edges[ (part1,part2) ]['evidList'].append( ev )
281                     edges[ (part1,part2) ]['intType'].append( ev['intType'] )
282                 else:
283                     edges[ (part1,part2) ] = {'evidList': [ ev ],
284                                                 'intType': [ev['intType']] }
285
286
287
288
289
290
291
```

## Network/Graph



# MIF File Parsing

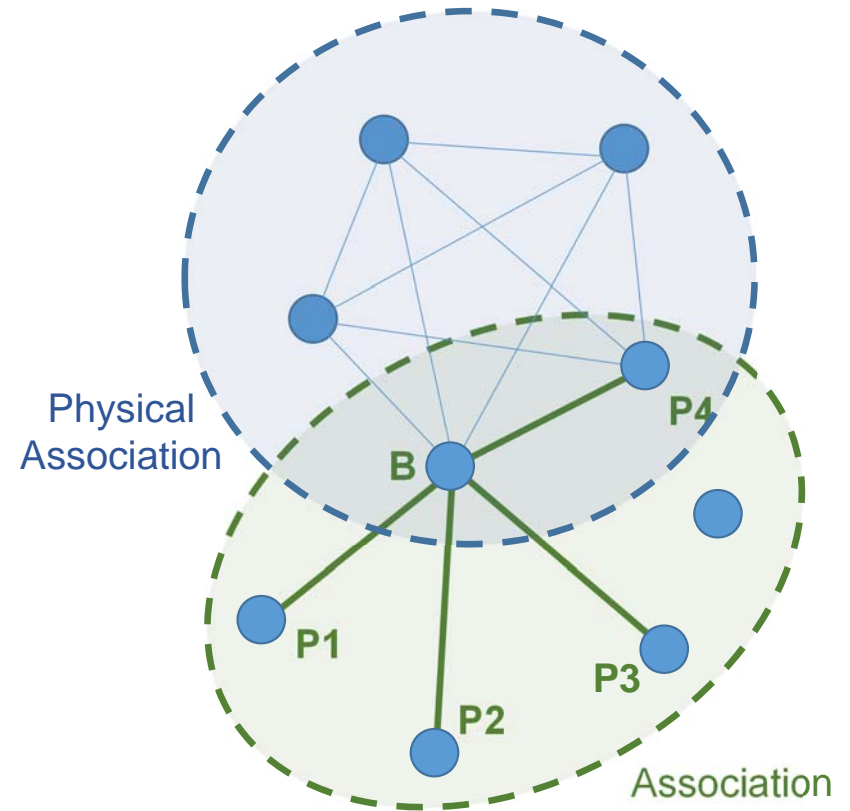
## Merging Interaction Data

```
333
334 def toJson( mifData, indent=None ):
335     return json.dumps( mifData, sort_keys=True, indent=indent )
336
337
338 def toSif( mifData, network ):
339     buffer = ''
340     for edge in network:
341         (p1,p2) = edge
342
343         p1type = mifData['irByAccId'][p1]['type']['name'][0]
344         p2type = mifData['irByAccId'][p2]['type']['name'][0]
345
346         buffer += "%s %s%s %s\n" % ( p1, p1type, p2type, p2 )
347     return buffer
348
349
350 def toEdgePoperties( mifData, network ):
351     buffer = '#acc\ttype id\ttype name\tmethod id\tmethod name\tpmid\n'
352     for edge in network:
353         (p1,p2) = edge
354         evList = network[edge]['evidList']
355         itypeList = network[edge]['intType']
356
357         p1type = mifData['irByAccId'][p1]['type']['name'][0]
358         p2type = mifData['irByAccId'][p2]['type']['name'][0]
359
360         line = "%s (%s%s) %s\t" % ( p1, p1type, p2type, p2 )
361
362         # type
363         typeId=''
364         typeName=''
365         for itype in itypeList:
366             typeId += itype['id'] + '|'
367             typeName += itype['name'] + '|'
368
369         line += typeId[:-1] + "\t" + typeName[:-1] + "\t"
370
371         # method
372         methodId=''
373         methodName=''
374         for ev in evList:
375             methodId += ev['expt']['detMeth']['id'] + '|'
376             methodName += ev['expt']['detMeth']['name'] + '|'
377
378         line += methodId[:-1] + "\t" + methodName[:-1] + "\t"
379
380         # source
381         source = ''
```

JSON

SIF

## Network/Graph



# MIF File Parsing

## Merging Interaction Data

```

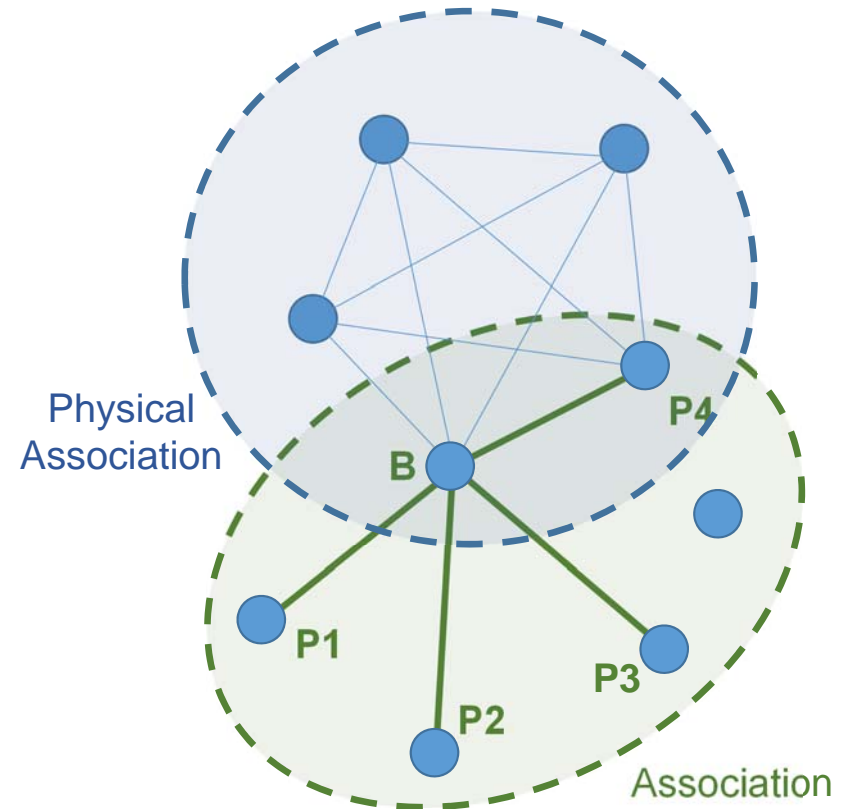
349
350 def toEdgeProperties( mifData, network ):
351     buffer = '#acc\ttype id\ttype name\tmethod id\tmethod name\tpmid\n'
352     for edge in network:
353         (p1,p2) = edge
354         evList = network[edge]['evList']
355         itypeList = network[edge]['intType']
356
357         p1type = mifData['irByAccId'][p1]['type']['name'][0]
358         p2type = mifData['irByAccId'][p2]['type']['name'][0]
359
360         line = "%s (%s%s) %s\t" % (p1, p1type, p2type, p2)
361
362         # type
363         typeId=''
364         typeName=''
365         for itype in itypeList:
366             typeId += itype['id'] + '|'
367             typeName += itype['name'] + '|'
368
369         line += typeId[:-1] + "\t" + typeName[:-1] + "\t"
370
371         # method
372         methodId=''
373         methodName=''
374         for ev in evList:
375             methodId += ev['expt']['detMeth']['id'] + '|'
376             methodName += ev['expt']['detMeth']['name'] + '|'
377
378         line += methodId[:-1] + "\t" + methodName[:-1] + "\t"
379
380         # source
381         source = ''
382         for ev in evList:
383             source += ev['expt']['pmid'] + '|'
384
385         line += source[:-1] + "\t"
386
387         buffer += line + "\n"
388
389     return buffer
390
391
392 def toNodeProperties( mifData, network ):
393     buffer = '#acc\tlabel\ttype\ttaxid\ttaxon label\ttaxon name\tgene\n'
394
395     nodeDict = mifData['irByAccId']
396
397     for node in nodeDict:

```

Edge Identifier

Edge Properties

## Network/Graph



# MIF File Parsing

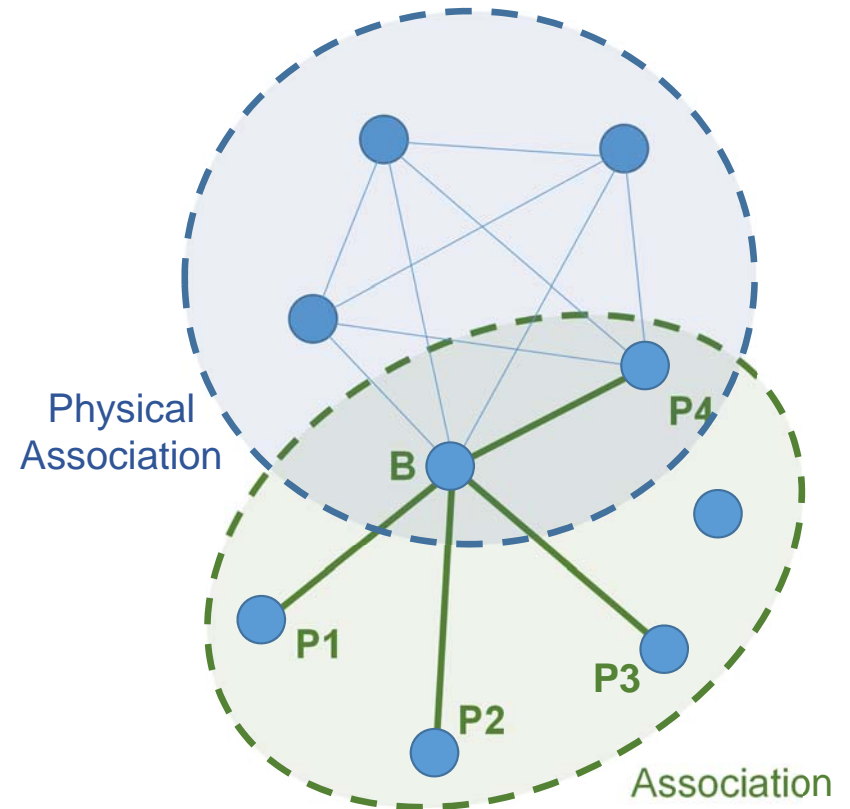
## Merging Interaction Data

```
391
392 def toNodeProperties( mifData, network ):
393     buffer = '#acc\tlabel\ttype\ttaxid\ttaxon label\ttaxon name\tgene\n'
394
395     nodeDict = mifData['irByAccId']
396
397     for node in nodeDict:
398         # accession
399         line = "%s\t" % ( nodeDict[node]['acc'] )
400
401         # label
402         line += "%s\t" % ( nodeDict[node]['label'] )
403
404         # type
405         line += "%s\t" % ( nodeDict[node]['type']['name'] )
406
407         # taxid
408         line += "%s\t" % ( nodeDict[node]['taxon']['taxid'] )
409
410         # taxon label
411         line += "%s\t" % ( nodeDict[node]['taxon']['label'] )
412
413         # taxon name
414         line += "%s\t" % ( nodeDict[node]['taxon']['name'] )
415
416         # gene
417         line += "%s\t" % ( nodeDict[node]['gene'] )
418
419         buffer += line + "\n"
420     return buffer
421
422
423 #-----
424 # merge files
425 #-----
426
427 mifData = {}
428
429 assert len(sys.argv) == 4
430 mydir = sys.argv[1]
431 mymth = sys.argv[2]
432 myout = sys.argv[3]
433
434 myfiles = []
435 for (dirpath, dirnames, filenames) in walk( mydir ):
436     myfiles.extend(filenames)
437     break
438
439 for f in myfiles:
```

Node Identifier

Node Properties

## Network/Graph



# Interaction Data Processing

*Take Home Message(s)*

- XML parsing is not that hard !!!
  - In many cases it is much easier than parsing text files
- Use example scripts as a starting point for your own
  - To read UniprotKB records
  - To read PSI-MI XML interaction files
  - To read any other XML (or HTML) files
- Be careful when dealing with
  - Potentially ambiguous and/or inconsistent protein identifiers
  - Interaction types
  - Multi-protein interactions (spoke/matrix expansion, changing interaction types)

# Interaction Data Processing

## *XML file parsing*

### Simple (and less simple but somewhat useful) project ideas

- List UniprotKB (gene names, Entrez gene identifiers, GO terms, etc) identifiers for each protein in a given MIF file
- Generate a FASTA file listing sequences all the bait proteins reported in a given MIF file
- Count interactions of a given protein (or a set of proteins) that are reported in a given MIF file
- Find all direct interactions of a given protein that are reported in a given MIF file
- List all proteins for each interaction reported in a given MIF file excluding proteins annotated (experimental role) as 'ancillary'
- Retrieve GO annotation from XML-formatted UniprotKB record – e.g.

<https://www.uniprot.org/uniprot/P60010.xml>

# Interaction Data Processing

## References

### XML Parsing

- lxml library: <https://lxml.de>
- XPath: <https://en.wikipedia.org/wiki/XPath> (and references therein)

### MIF (miXML) format specification

- XSD: <https://github.com/HUPO-PSI/miXML> (2.5/src, 3.0/src directories)
- Publications
  - Hermjakob H *et al.* The HUPO PSI's molecular interaction format--a community standard for the representation of protein interaction data.  
Nat Biotechnol. 22:177-83 (2004). PMID: 14755292
  - Kerrien S *et al.* Broadening the horizon--level 2.5 of the HUPO-PSI format for molecular interactions.  
BMC Biol. 5:44 (2007). PMID: 17925023
  - Sivade Dumousseau M *et al.* Encompassing new use cases - level 3.0 of the HUPO-PSI format for molecular interactions.  
BMC Bioinformatics. 19:134. doi: 10.1186/s12859-018-2118-1 (2018). PMID: 29642841