

Matúš Kalaš<sup>1</sup>, Sveinung Gundersen<sup>2</sup>, László Kaján<sup>3</sup>, Hervé Ménager<sup>4</sup>, Jon Ison<sup>5</sup>, Christophe Blanchet<sup>6</sup>, Steve Pettifer<sup>7</sup>, Rodrigo Lopez<sup>8</sup>, Kristoffer Rapacki<sup>5</sup>, Inge Jonassen<sup>1</sup>, and open for contributions

2017

<sup>1</sup>Computational Biology Unit, Department of Informatics, University of Bergen, Bergen, Norway; <sup>2</sup>Institute for Cancer Research, Oslo University Hospital and Department of Informatics, University of Oslo, Oslo, Norway; <sup>3</sup>unaffiliated, previously Bioinformatics and Computational Biology Department, Technische Universität München, Garching, Germany; <sup>4</sup>Institut Pasteur, Paris, France; <sup>5</sup>Department of Systems Biology, Technical University of Denmark, Kongens Lyngby, Denmark; <sup>6</sup>L'Institut Français de Bioinformatique, Gif-sur-Yvette, France; <sup>7</sup>School of Computer Science, University of Manchester, Manchester, UK; <sup>8</sup>European Bioinformatics Institute, EMBL, Wellcome Trust Genome Campus, Hinxton, Cambridge, UK.

[/bioxsd/bioxsd](https://github.com/bioxsd/bioxsd)

[@BioXSD](https://twitter.com/BioXSD)

<http://groups.google.com/group/bioxsd>

<http://bioxsd.org>

[support@bioxsd.org](mailto:support@bioxsd.org)

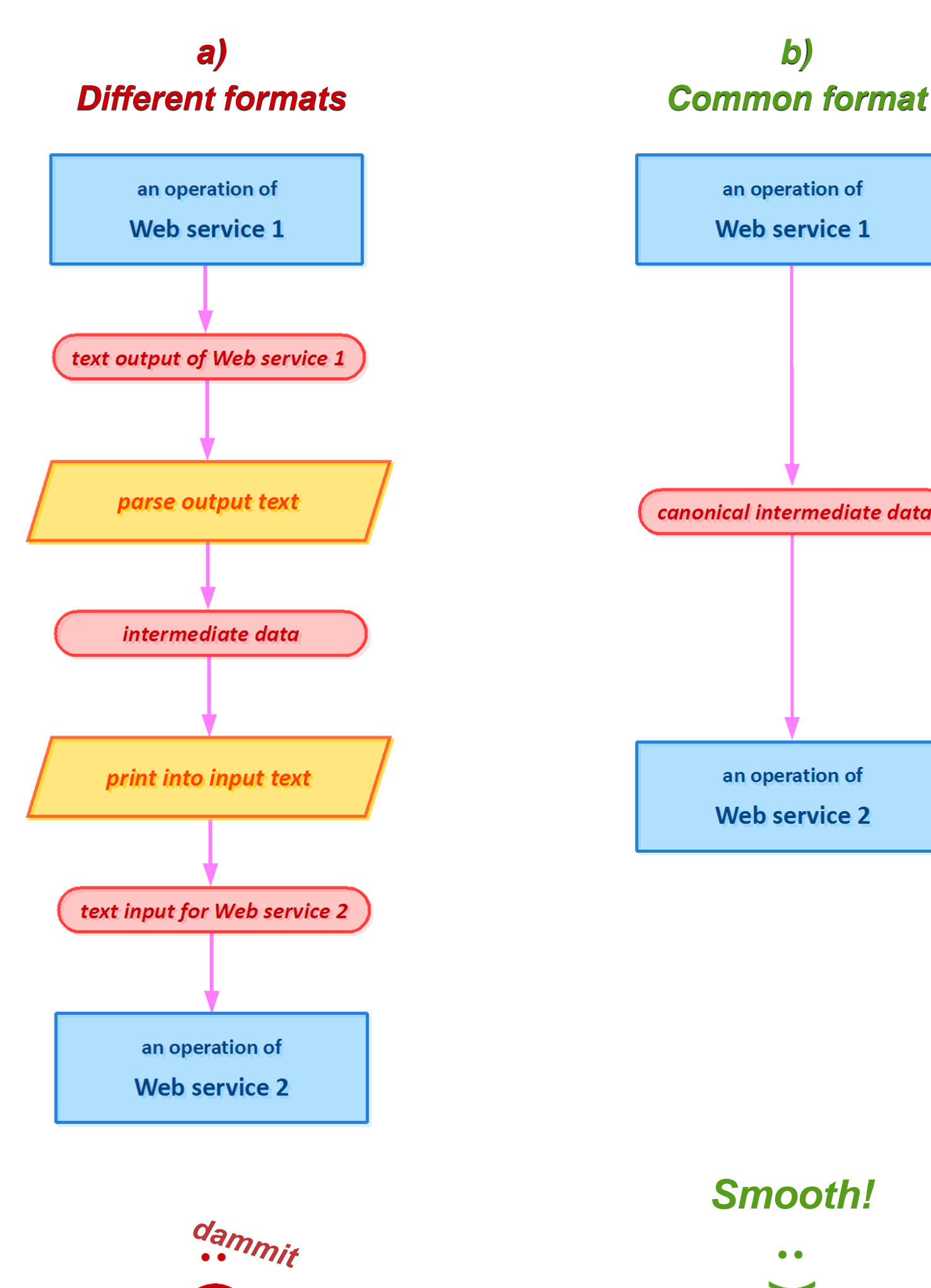
Latest stable release: <http://bioxsd.org/BioXSD-1.1.xsd>

## MOTIVATION

Without a common format, using diverse tools in a workflow demands conversions, “shims”, or do-it-yourself parsing. And worst of all, maintaining these in the future.

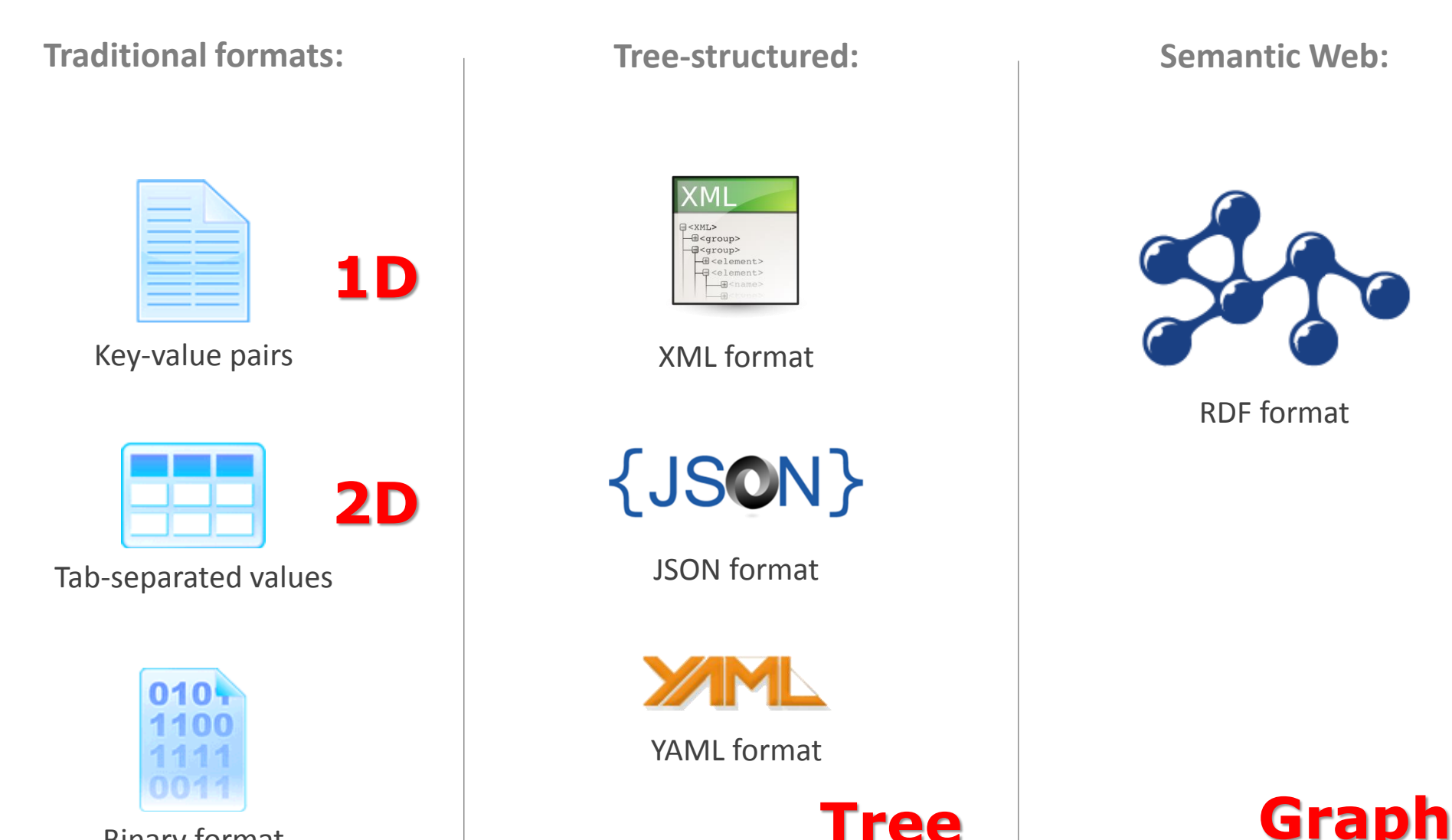
The 2 scenarios show demands for connecting 2 tools (e.g. Web services) that use:

- a) Different formats
- b) A common format



## TECHNOLOGY CHOICES

Different paradigms of data formatting represent data differently.



A machine-understandable definition of a specific format (a data model, a schema) is highly beneficial for validation and maintainability.

<b>GTrack format</b> TSV with column definitions <a href="http://gtrack.no">http://gtrack.no</a>	<b>XML Schema (XSD) 1.0</b> <b>XML Schema 1.1</b> <b>Relax NG</b> <b>JSON Schema</b> ...	<b>OWL</b> ...
--	--	-------------------

## SIMPLE EXAMPLE: BioXSD sequence record

Example data instance, **BioXSD** in XML:

```
<mySequenceRecord
  xsi:type="bc:GeneralAminoacidSequenceRecord"
  xmlns:bc="http://bioxsd.org/BioXSD-1.1"
  xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xsi:schemaLocation="http://bioxsd.org/BioXSD-1.1 http://bioxsd.org/BioXSD-1.1.xsd"
>
  <bc:sequence>MDPLGDTLRLRLREAFHAGRTRPAEFRAAQGLGRFLQENKQLLHDAL</bc:sequence>
  <bc:species
    dbName="NCBI Taxonomy"
    accession="9606"
    entryUri="http://www.ncbi.nlm.nih.gov/Taxonomy/Browser/wwwtax.cgi?id=9606"
    speciesName="Human"
  />
  <bc:reference
    dbName="UniProt"
    accession="P43353"
    entryUri="http://www.uniprot.org/uniprot/P43353"
    sequenceVersion="1"
    variantAccession="P43353-1"
  />
  <bc:subsequencePosition
    bc:segment min="1" max="48"/>
  </bc:subsequencePosition>
  <bc:reference
    bc:name="Aldehyde dehydrogenase family 3 member B1 (ALDH3B1), N-terminus" bc:name>
</mySequenceRecord>
```

In **BioJSON**:

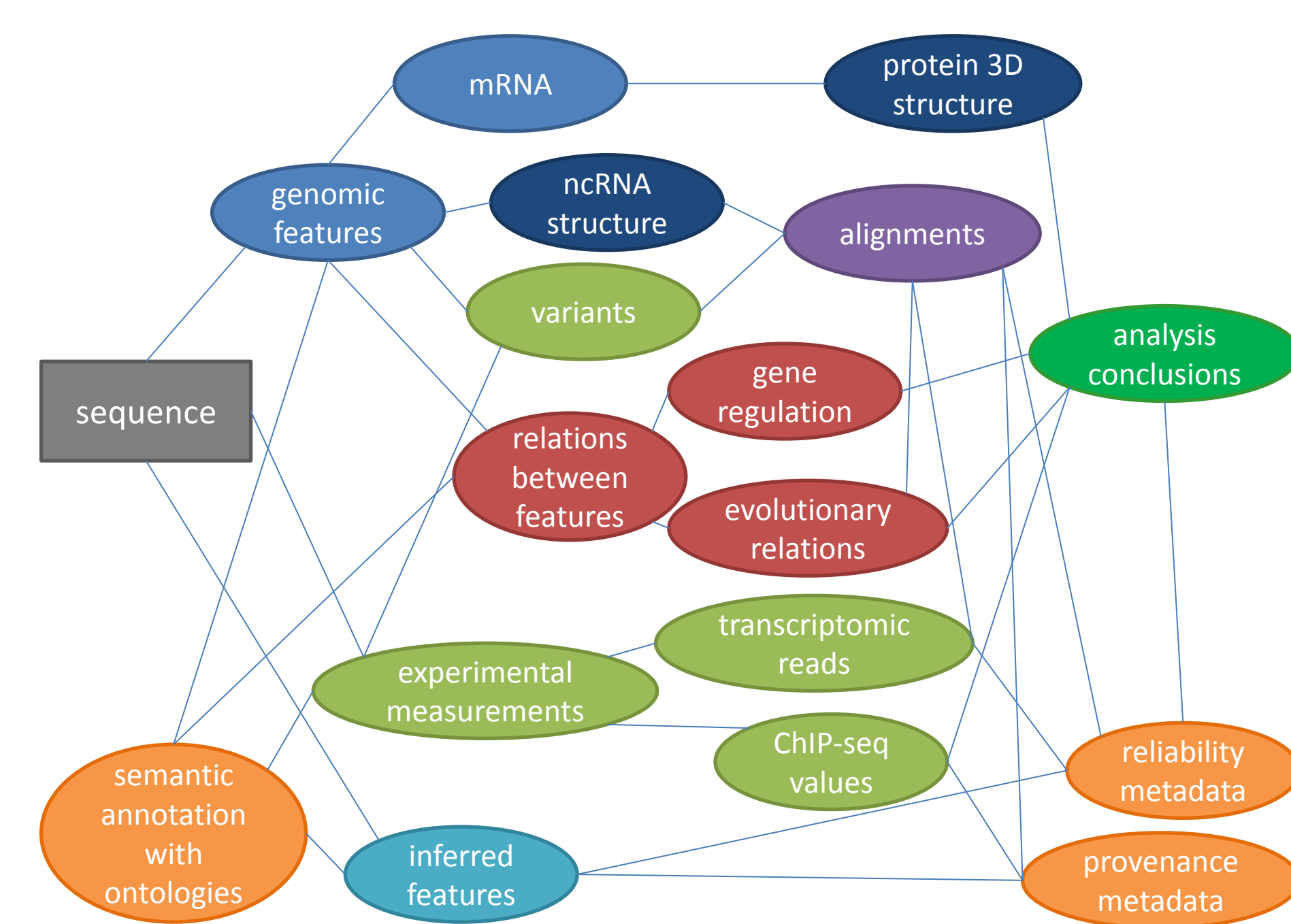
```
{
  "sequence": "MDPLGDTLRLRLREAFHAGRTRPAEFRAAQGLGRFLQENKQLLHDAL",
  "species": {
    "dbName": "NCBI Taxonomy",
    "accession": "9606",
    "entryUri": "http://www.ncbi.nlm.nih.gov/Taxonomy/Browser/wwwtax.cgi?id=9606",
    "speciesName": "Human"
  },
  "reference": {
    "dbName": "UniProt",
    "accession": "P43353",
    "entryUri": "http://www.uniprot.org/uniprot/P43353",
    "sequenceVersion": "1",
    "variantAccession": "P43353-1",
    "subsequencePosition": {
      "segment": {
        "min": 1,
        "max": 48
      }
    }
  },
  "name": "Aldehyde dehydrogenase family 3 member B1 (ALDH3B1), N-terminus"
}
```

In **BioYAML**:

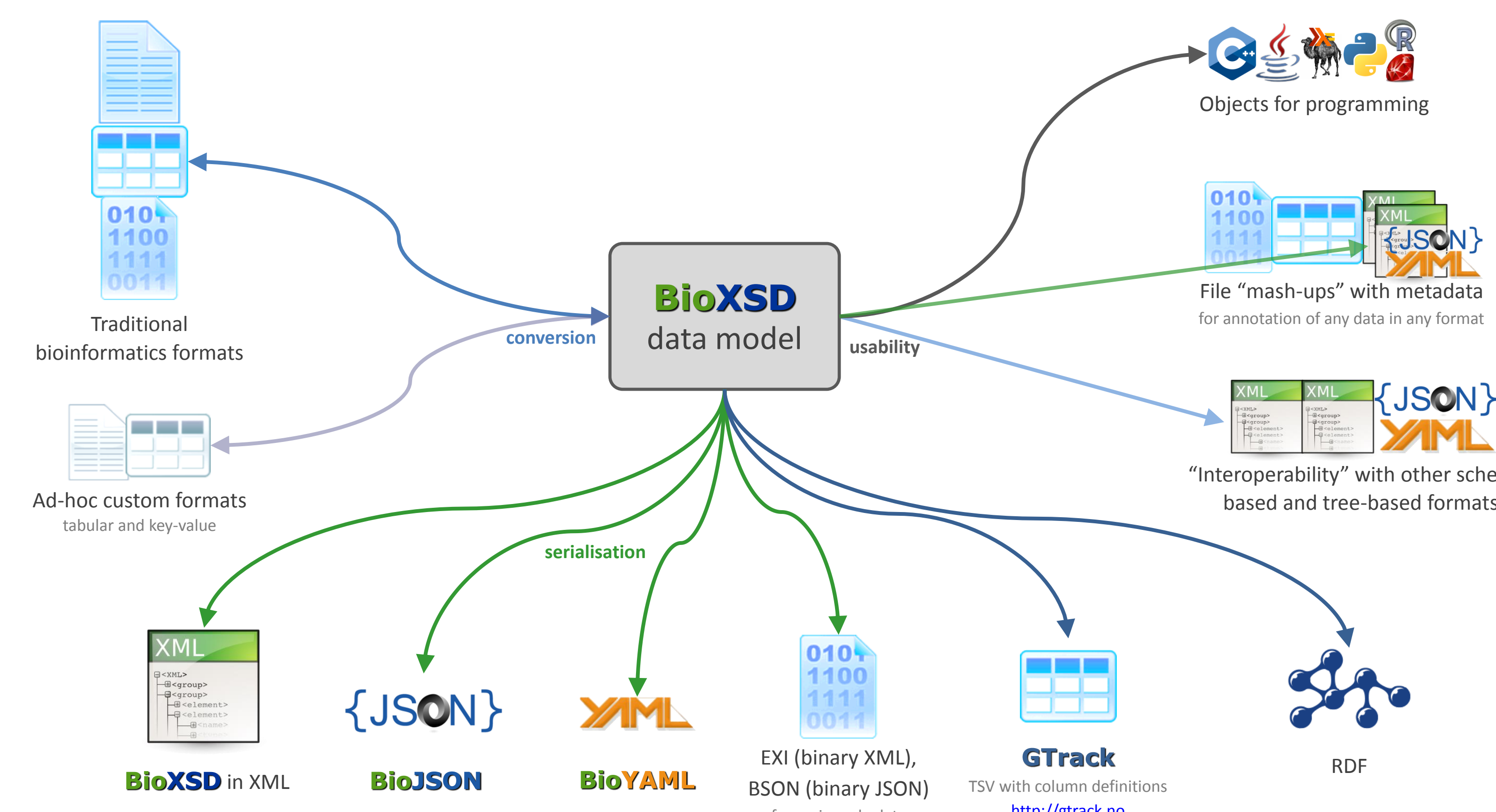
```
---
sequence: MDPLGDTLRLRLREAFHAGRTRPAEFRAAQGLGRFLQENKQLLHDAL
species:
  dbName: NCBI Taxonomy
  accession: 9606
  entryUri: http://www.ncbi.nlm.nih.gov/Taxonomy/Browser/wwwtax.cgi?id=9606
  speciesName: Human
reference:
  dbName: UniProt
  accession: P43353
  entryUri: http://www.uniprot.org/uniprot/P43353
  sequenceVersion: 1
  variantAccession: P43353-1
  subsequencePosition:
    segment:
      min: 1
      max: 48
name: Aldehyde dehydrogenase family 3 member B1 (ALDH3B1), N-terminus
```

## COMPLEX EXAMPLE: BioXSD feature record

BioXSD can represent diverse interconnected features of a sequence, together with related data and metadata, in an integrated data record.



## ONGOING DEVELOPMENTS: single data model with multiple choices of exchange formats and conversions



**BioXSD** has been developed as a tree-based data model and an exchange format for basic bioinformatics data, centred around a bio-polymer sequence. BioXSD allows integration of diverse features, information, measurements, and inferred values about a biological molecule or its part, annotated with provenance and reliability metadata, ontology concepts, scientific remarks, and conclusions.

BioJSON and BioYAML are the ongoing developments. These exchange formats are based on the same data model as BioXSD, but providing serialisations in JSON and YAML respectively. BioJSON and BioYAML thus enrich the BioXSD family with alternatives to the original XML.

As tree-based data formats, BioXSD, BioJSON, and BioYAML are particularly suitable for programming in object-oriented languages, and for use with web applications and web APIs (Web services), while at the same time allowing a reasonable level of human readability.

BioXSD|BioJSON|BioYAML are developed together with **GTrack** (the universal tabular format for sequence features). The **BioXSD|GTrack family**, as a “service” developed by **ELIXIR Norway**, is going to support smooth interoperability between these alternative, universal formats, and between the tools that consume or provide them as inputs or outputs.