# UPOV

International Union for the Protection of New Varieties of Plants

| | |
|---|---|
| **Working Group on Biochemical and Molecular Techniques and DNA-Profiling in Particular** | **BMT/16/5 Add.** |

| | |
|---|---|
| **Sixteenth Session** <br> **La Rochelle, France, November 7 to 10, 2017** | **Original:** English <br> **Date:** November 3, 2017 |

**ADDENDUM TO**
**STANDARDS FOR DATABASES CONTAINING MOLECULAR INFORMATION**

*prepared by the Office of the Union*

*Disclaimer: this document does not represent UPOV policies or guidance*

The Annex of this document contains a copy of a presentation "Standards for databases containing molecular information" to be made by the Office of the Union at the sixteenth session of the Working Group on Biochemical and Molecular Techniques, and DNA-Profiling in Particular (BMT).

[Annex follows]

# Standards for databases containing molecular information

November 7, 2017

**UPOV**
International Union for the Protection of New Varieties of Plants

---

## PREVIEW

- Databases
- WIPO ST.26
- WIPO ST.26 Software

2

# Databases

- Organized array of information
- Place where you put things in, and you should be able to get them out again.
- Allows you to search.

UPOV

3

# What you can store

- Fingerprints
  - 1-D electrophoresis gels scanned as bitmaps (RFLP, PFGE, Ribotyping, RAPD, DGGE & TGGE, etc.)
  - Sequencer chromatogram files (AFLP, VNTR, HDA, etc.)
  - Spectrophotometric files
  - MALDI & SELDI profiles
  - All other kinds of densitometric profiles

- Character data : Phenotypic test panels
  - Antibiotic resistance profiles
  - Fatty acid and quinolone profiles
  - Hybridization blots
  - Biochemical & morphological features
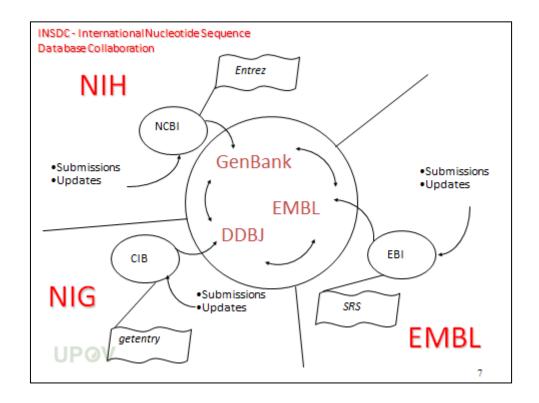  - Microarray & Genechip data

UPOV

# What you can store (cont'd)

- Sequence data
  - Sequence trace (chromatogram) files
  - Formatted sequences from public databases (EMBL, GenBank)
  - Aligned sequences
  - Amino acid sequences

UPOV

# Database Examples in Bioinformatics

|  | Primary database | Secondary database |
|---|---|---|
| Synonyms | Archival database | Curated database; knowledgebase |
| Source of data | Direct submission of experimentally-derived data from researchers | Results of analysis, literature research and interpretation, often of data in primary databases |
| Examples | ✓ GenBank/EMBL/DDBJ (nucleotide sequence)<br>✓ Protein Data Bank (PDB, coordinates of three-dimensional macromolecular structures)<br>✓ Medline (literature)<br>✓ IMEx databases (protein interactions)<br>✓ ArrayExpress Archive and GEO (functional genomics data) | ✓ InterPro(protein families, motifs and domains)<br>✓ UniProt Knowledgebase - SwissProt (sequence and functional information on proteins)<br>✓ Ensembl (variation, function, regulation and more layered onto whole genome sequences) |

UPOV

INSDC - International Nucleotide Sequence Database Collaboration

NIH

Entrez

NCBI

GenBank

EMBL

DDBJ

•Submissions
•Updates

•Submissions
•Updates

CIB

EBI

NIG

•Submissions
•Updates

SRS

UPOV

getentry

EMBL

7



## PREVIEW

- Databases
- WIPO ST.26
- WIPO ST.26 Software

8

## What is WIPO ST.26?

- ST.26 is the recommended standard for the presentation of nucleotide and amino acid sequence listings using XML

- It defines the sequence disclosures in a patent application required to be included in a sequence listing

9

## WIPO ST.26

- Based on INSDC specifications
- Faciliates searching of the sequence data
- Allows sequence data to be exchanged in electronic form and introduced into computerized databases.

UPOV

## Sequence Listing in XML
### General information part

- ApplicationIdentification : Mandatory
  - IPOfficeCode
  - ApplicationNumberText
  - FilingDate
- ApplicantFileReference: Optional
- EarliestPriorityApplicationIdentification : Mandatory if Priority is claimed
- ApplicantName : Mandatory
- ApplicantNameLatin : Optional
- InventorName: Optional
- InventorNameLatin: Optional
- InventionTitle: Mandatory in the language of filing
- SequenceTotalQuantity: Mandatory

UPOV

## Sequence Listing in XML
### Sequence Data part

- One or more SequenceData elements
- Each SequenceData has a mandatory attribute sequenceIDNumber

| Element | Description | Mandatory/Not Included | |
|---|---|---|---|
| | | Sequences | Intentionally Skipped Sequences |
| INSDSeq_length | Length of the sequence | Mandatory | Mandatory with no value |
| INSDSeq_moltype | Molecule type | Mandatory | Mandatory with no value |
| INSDSeq_division | Indication that a sequence is related to a patent application | Mandatory with the value "PAT" | Mandatory with no value |
| INSDSeq_feature-table | List of annotations of the sequence | Mandatory | Must NOT be included |
| INSDSeq_sequence | Sequence | Mandatory | Mandatory with the value "000" |

UPOV

# Feature Keys and Qualifiers

- Nucleic Acid Sequences
  - Agreed upon by the International Nucleotide Sequence Database Collaboration (INSDC)
  - 49 feature keys and 80 qualifiers for nucleic acid sequences: INSDC feature keys/qualifiers not relevant for patent data not included

UPOV

June 2012

13

# Sequence Listing in XML
## Sequence Data part
### Feature Table

- Information on location and roles of various regions within a particular sequence
- One or more INSDFeature elements

| Element | Description | Mandatory/Optional |
|---|---|---|
| INSDFeature_key | A word or abbreviation indicating a feature | Mandatory |
| INSDFeature_location | Region of the presented sequence which corresponds to the feature | Mandatory |
| INSDFeature_quals | Qualifier containing auxiliary information about a feature | Mandatory where the feature key requires one or more qualifiers, e.g. source; otherwise, Optional |

UPOV

# Variety

| | |
|---|---|
| qualifier | variety |
| definition | variety (= varietas, a formal Linnaean rank) of organism from which sequence was derived. |
| value format | free text<br>(NOTE: this value may require translation for national/regional procedures) |
| example | <INSDqualifier_value>insularis</INSDqualifier_value> |
| comment | use the cultivar qualifier for cultivated plant varieties, i.e., products of artificial selection; varieties other than plant and fungal varietas should be annotated via a note qualifier, e.g. with the value<br><INSDqualifier_value>breed:cukorova</INSDqualifier_value> |

UPOV

---

# Example: PP28388

- Variety: CIMAP-KHUSINOLIKA
- Species/Crop: VETIVER ( CHRYSOPOGON ZIZANIODES )
- Phenotype: PRODUCES KHUSINOL RICH

ESSENTIAL OIL UNDER SHORT DURATION CULTIVATION

- What is stored: ISSR-PCR primers

UPOV

## ISSR Primer

```
<INSDQualifier>
    <INSDQualifier_name>note</INSDQualifier_name>
    <INSDQualifier_value>A synthetic ISSR Primer</INSDQualifier_value>
</INSDQualifier>
</INSDFeature_quals>
</INSDFeature>
</INSDSeq_feature-table>
<INSDSeq_sequence>agagagagag agagagt </INSDSeq_sequence>
</INSDSeq>
</SequenceData>
```

# Example: PP16174

- Variety: B12

- Sepecies/Crop: ST. AUGUSTINE GRASS

- Prior application number: AU PBR 2002/342

- What is stored: Primer
  - ccgcatctac

UPOV

## Primer

```
                    <INSDQualifier>
                        <INSDQualifier_name>note</INSDQualifier_name>
                        <INSDQualifier_value>Primer</INSDQualifier_value>
                    </INSDQualifier>
                </INSDFeature_quals>
            </INSDFeature>
        </INSDSeq_feature-table>
        <INSDSeq_sequence>ccgcatctac</INSDSeq_sequence>
    </INSDSeq>
</SequenceData>
</ST26SequenceListing>
```

UPOV

## Example: PP15792

- Variety: BEINEKE 8
- Species/Crop: Black walnut
- What is stored: 18 primers
    - gacgacgaag gtgtacggat
    - ccatgaaact tcatgcgttg
    - ......
    - ttgaacaaaa ggccgttttc

UPOV

# Example: PCT/US2015/055339

- TOMATO PLANTS WITH IMPROVED DISEASE RESISTANCE
- UPOV TG/44/11 Char 57: resistance to Tomato yellow leaf curl virus
- What is stored: probes and primers

UPOV

## PREVIEW

- Databases
- WIPO ST.26
- WIPO ST.26 Software

**UPOV**

25

---

# WIPO ST.26 Software

- Editing or importing sequences in ST.26 format
- Validation of sequences
- Transformation of ST.25 sequences to ST.26
- Importing existing sequence data in industry format, e.g. GenBank, EMBL and FASTA
- Presentation of XML in human readable format
- Multi language support: interface, message
- "Free text" translation support (the "free text" must be in Basic Latin in the sequence listing)

**UPOV**

# Timelines

- End of 2017: Proof of concept
- 2018: Testing and upgrades

UPOV

[End of Annex and of document]