

# Widening the impact of DOI for traceability and ABS

## Global Policy Meets Data Interoperability



Catalysing advanced conservation, management and traceability of PGR  
DivSeek Workshop | January 18, 2022



**Pankaj Jaiswal**



**Daniele Manzella  
Marco Marsella**



**Elizabeth Arnaud  
Brian King**



# Global Information System of the ITPGRFA (GLIS)

- Launched Oct. 2017
- Over 1.1M DOIs assigned
- All of CGIAR
- Research centres and Universities
- National collections
- Started with *ex situ* and breeding
- Soon *in situ* and on farm too
- DOIs to address identification
- Link to web resources
- Link to publications and datasets

The screenshot shows the PGRFA DOI page for 10.18730/5ER3F. The page is titled "PGRFA doi:10.18730/5ER3F" and includes a citation: "Citation: https://doi.org/10.18730/5ER3F". The page is organized into several sections:

- Main descriptors:** Breeding, DOI info
- Organization/individual conserving the PGRFA:** International Rice Research Institute, DAPO BOX 7777, 1301 Metro Manila, Philippines. WIEWS code: PHL001 [Details], Easy-SMTA PID: 00AB40.
- Local identifier:** IRGC 127122
- Date:** 2011-05-01
- Creation method:** In-house variant from 10.18730/3F11~
- Taxon:** Oryza sativa Linnaeus
- Common name:** Rice
- Biological status:** Genetic stock
- Names:** ANADI WHITE::IRGC 61897-1
- Other identifiers:**
  - MLS status:** Art. 15 collection
  - Historical:** No

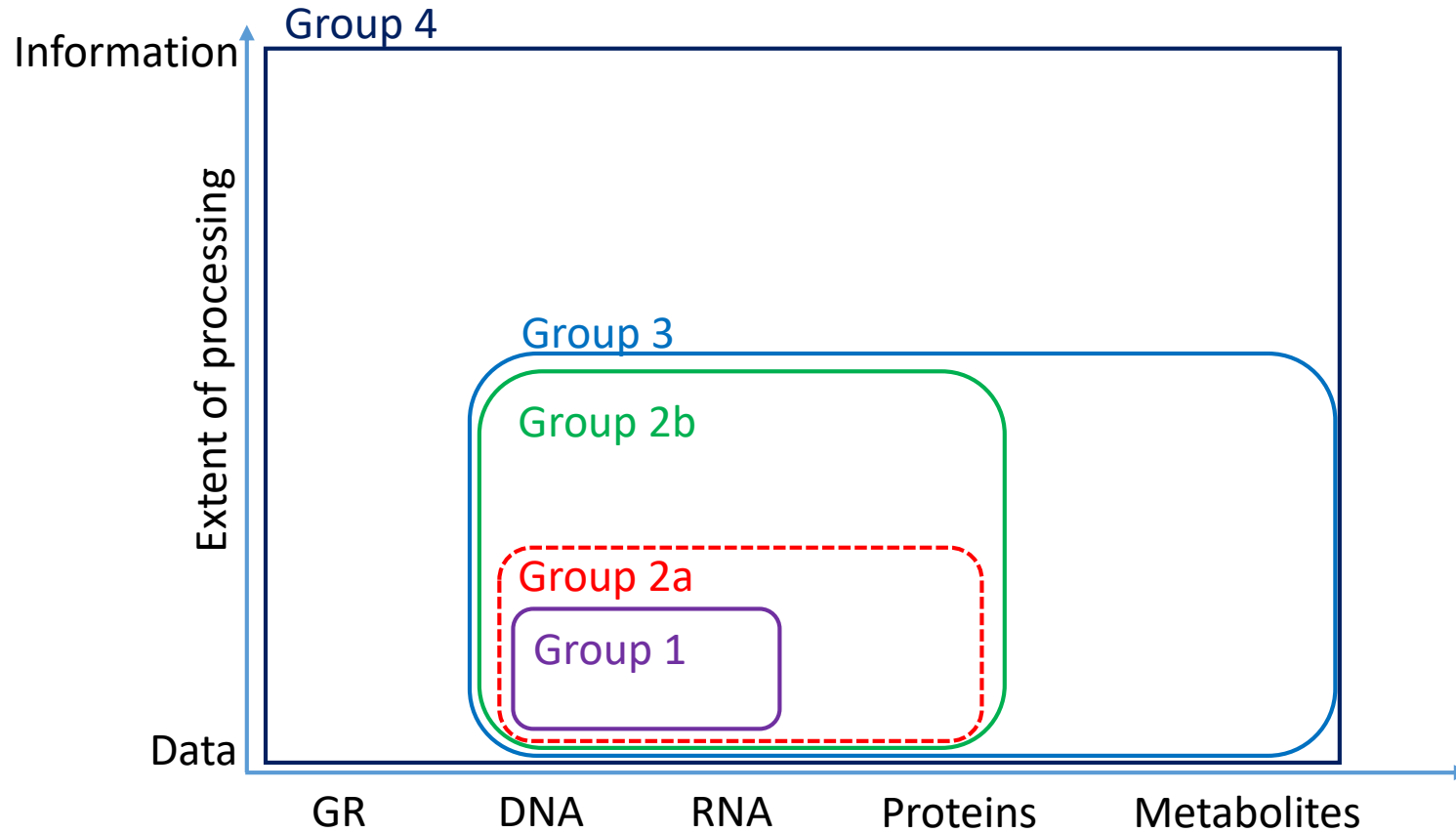
Below the main descriptors, there is a section for "Links to associated information (1-3 of 3)":

Keywords	URL
Passport data	<a href="http://www.fao.org/wiews/data/ex-situ-sdg-251/search/en/?doi=10.18730/5ER3F#details">http://www.fao.org/wiews/data/ex-situ-sdg-251/search/en/?doi=10.18730/5ER3F#details</a>
Passport data	<a href="https://www.genesys-pgr.org/10.18730/5ER3F">https://www.genesys-pgr.org/10.18730/5ER3F</a>
Genomics	<a href="https://snp-seek.irri.org/_variety.zul?irrsid=313-11624">https://snp-seek.irri.org/_variety.zul?irrsid=313-11624</a>

Below the links, there is a section for "Publications and datasets citing this PGRFA (1-6 of 6)":

Type	Title	Published	Journal	Authors	Publisher
Paper	<a href="#">Variation in seed longevity among diverse Indica rice varieties</a>	2019-06-10	Annals of Botany	Jae-Sung Lee, Marlina Velasco-Punzalan, Myrish Pacleb, Rocel Valdez, Tobias Kretzschmar, Kenneth L McNally, Abdel M Ismail, Pompe C Sta. Cruz, N Ruairaidh Sackville Hamilton, Fiona R Hay	Oxford Academic
Paper	<a href="#">An imputation platform to enhance integration of rice genetic resources</a>	2018-08-25	Nature Communications	Diane R. Wang, Francisco J. Agosto-Pérez, Dmytro Chebotarov, Yuxin Shi, Jonathan Marchini, Melissa Fitzgerald, Kenneth L. McNally, Nikolai Alexandrov, Susan R. McCouch	Nature Research
Paper	<a href="#">Seed longevity phenotyping: recommendations on research methodology</a>	2018-05-11	Journal of Experimental Botany	Fiona R. Hay, Rocel Valdez, Jae-Sung Lee, Pompe C. Sta. Cruz	Society for Experimental Biology
Paper	<a href="#">The 3,000 rice genomes project: new opportunities and challenges for future rice research</a>	2014-05-28	GigaScience	Jia-Yang Li, Jun Wang, Robert S Ziegler	Oxford University Press
Paper	<a href="#">The 3,000 rice genomes project</a>	2014-05-28	GigaScience	CAAS, BGI, IRRI	Oxford University Press
Dataset	<a href="#">The Rice 3000 Genomes Project Data</a>	2014-05-27			GigaScience Database

# Digital Sequence Information



Adapted from Houssen et al. (2020)

## Granular Options for Subject Matter Groupings

**Group 1** - Narrow: DNA and RNA

**Group 2a** includes DNA/RNA sequence data including non-coding sequences, and information on the sequence assembly, including structural annotation and genetic mapping, as well as protein sequence data.

**Group 2b** is the same as group 2a in addition to which it includes functional annotation of genes, gene expression information, epigenetic data, and molecular structures of proteins.

**Group 3** is the same as group 2b, but adds data on other macromolecules and metabolites, including their molecular structures

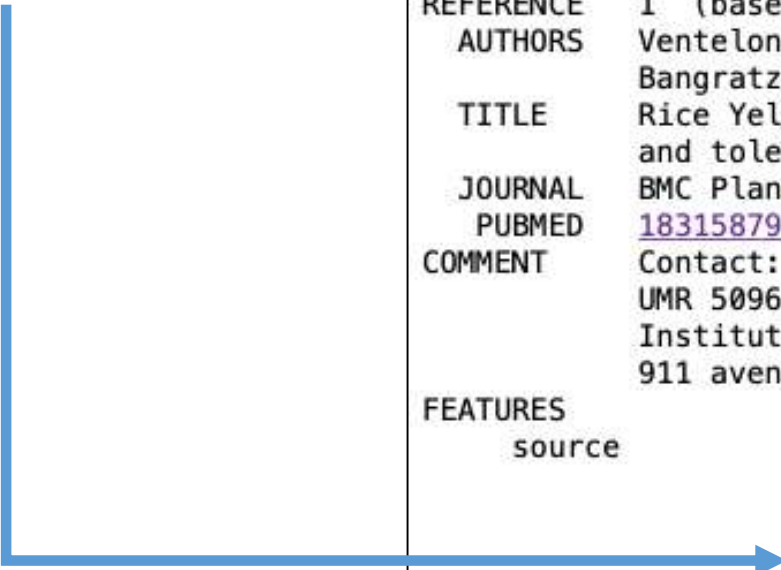
**Group 4** – Broad: DNA, RNA, protein, metabolites + traditional knowledge, ecological interactions, etc.





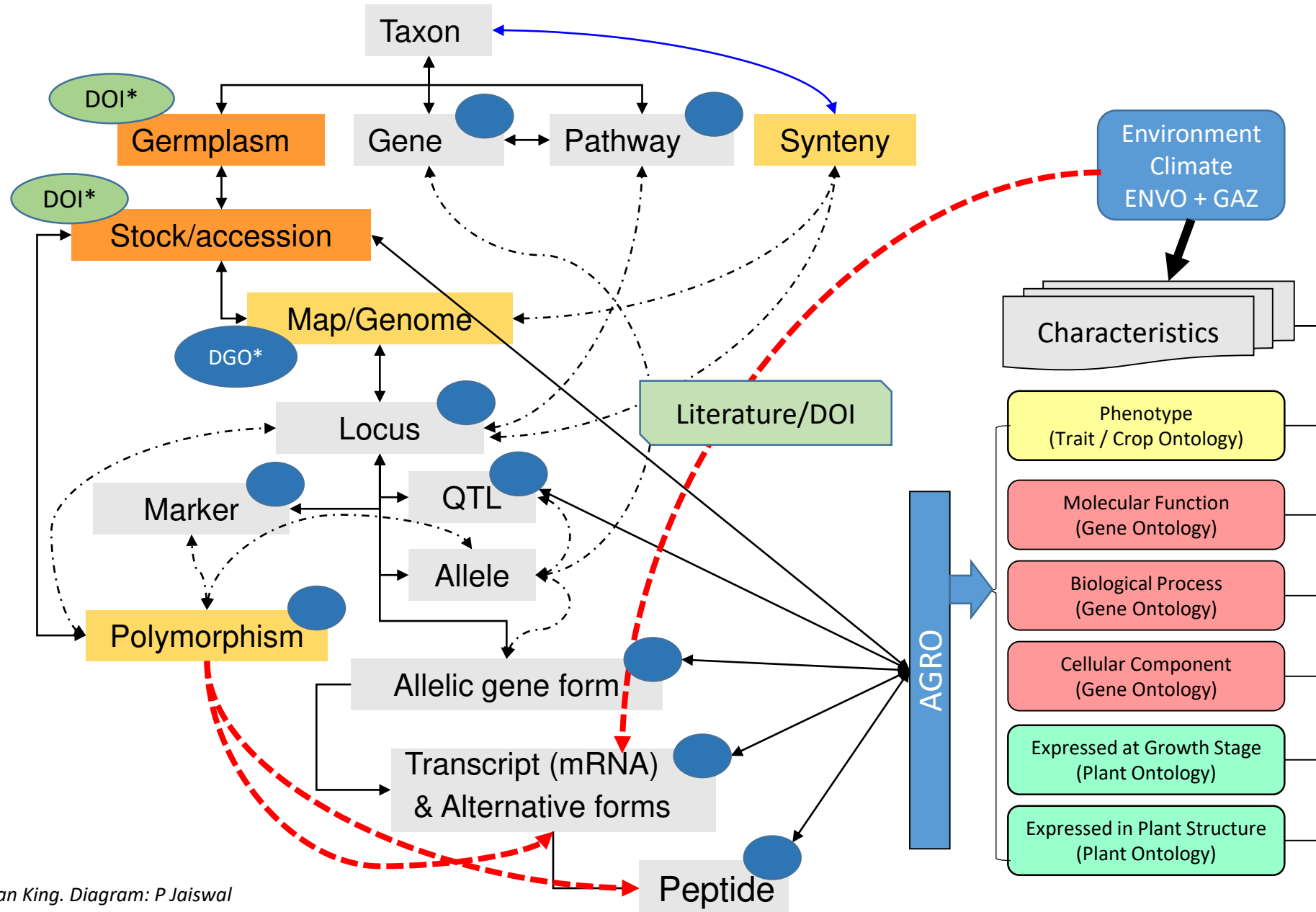
# International Nucleotide Sequence Database (INSDC) Collaboration

Source material is referenced through this attribute which is not mandatory



```
LOCUS      DQ884074                257 bp    mRNA     linear   EST 24-FEB-2011
DEFINITION DQ884074 Oryza sativa (indica cultivar-group) cv. IR64 cDNA-AFLP
           fragment Oryza sativa Indica Group cDNA clone 51_9b, mRNA sequence.
ACCESSION  DQ884074
VERSION    DQ884074.1
DBLINK     BioSample: SAMN00165158
KEYWORDS   EST.
SOURCE     Oryza sativa Indica Group (long-grained rice)
           ORGANISM Oryza sativa Indica Group
           Eukaryota; Viridiplantae; Streptophyta; Embryophyta; Tracheophyta;
           Spermatophyta; Magnoliopsida; Liliopsida; Poales; Poaceae; BOP
           clade; Oryzoideae; Oryzeae; Oryzinae; Oryza; Oryza sativa.
REFERENCE  1 (bases 1 to 257)
AUTHORS    Ventelon-Debout,M., Tranchant-Dubreuil,C., Nguyen,T.-T.-H.,
           Bangratz,M., Sire,C., Delseny,M. and Brugidou,C.
TITLE      Rice Yellow Mottle Virus stress responsive genes from susceptible
           and tolerant rice genotypes
JOURNAL    BMC Plant Biol. 8, 26 (2008)
PUBMED     18315879
COMMENT    Contact: Ventelon-Debout M
           UMR 5096
           Institut de Recherche pour le Developpement
           911 avenue d'Agropolis, BP54501, Montpellier, 34394, France.
FEATURES   source
           Location/Qualifiers
           1..257
           /organism="Oryza sativa Indica Group"
           /mol_type="mRNA"
           /cultivar="IR64"
           /db_xref="taxon:39946"
           /clone="51_9b"
           /clone_lib="SAMN00165158 Oryza sativa (indica
           cultivar-group) cv. IR64 cDNA-AFLP fragment"
```

# What we propose to build



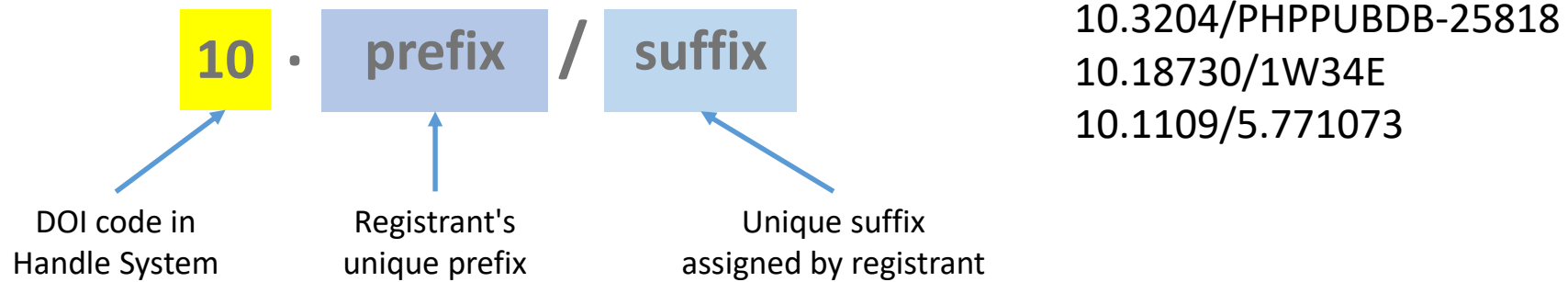
# What is a Digital Genomic Object (DGO)?

---

A DSI element is a DGO identified by a DOI and collected according to standard/defined procedure with very defined types of component data, can be linked via semantics

# The DOI System

DOIs (Digital Object Identifiers) are a kind of Permanent Unique identifier (PID)



The International DOI Foundation maintains the infrastructure

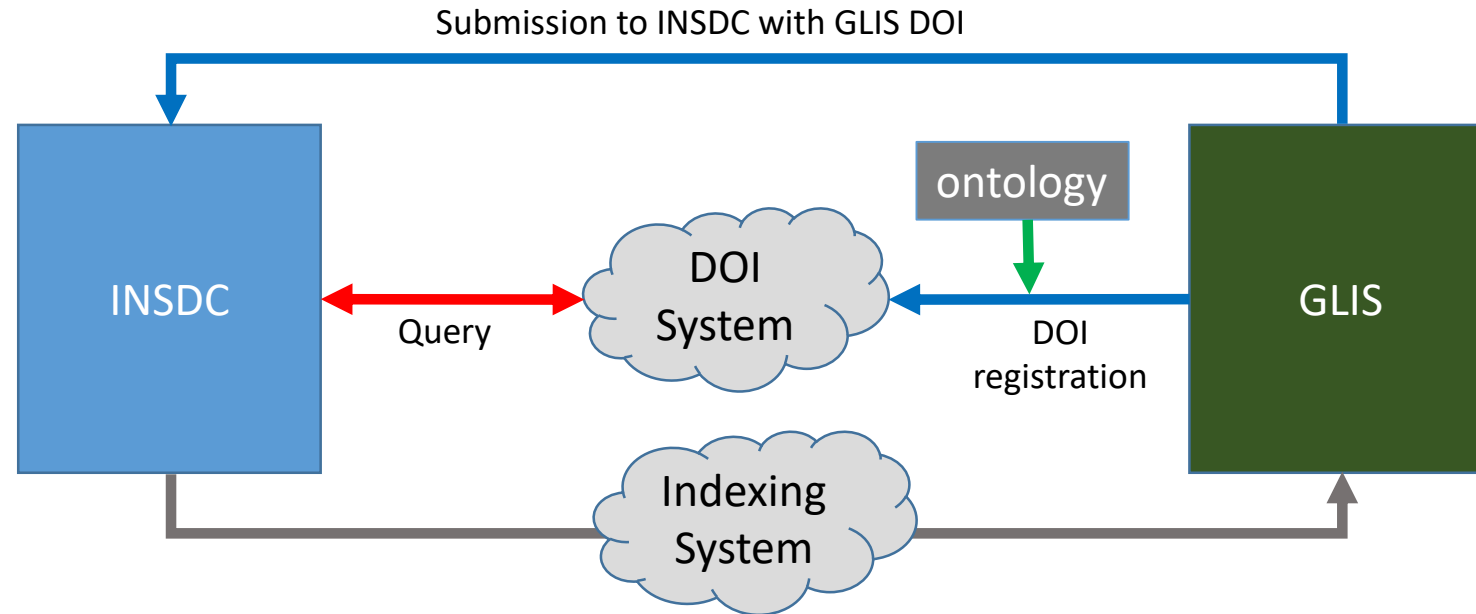
DOIs are assigned by Registration Agencies such as Crossref or DataCite

The DOI System offers concepts and tools suitable to our needs

- Flexible metadata structure to describe the different DSI components
- Provisions for establishing relationships among DOIs and other PIDs
- PID Graph, funded by EU project FREYA, records references among PIDs
- Make Data Count, records citation metrics for datasets
- DOIs are already being assigned to
  - Plant Genetic Resources
  - Papers and publications
  - Datasets

# Our pilot proposal

Scope: Select a couple of species data sources



- Researcher obtains DOI for source material from GLIS is not already available
- Researcher submits dataset to INSDC referencing GLIS DOI for source material
- GLIS discovers INSDC records referencing its own DOIs using an Indexing System
- If INSDC opts to assign DOIs to its Accessions, the Indexing System is not required



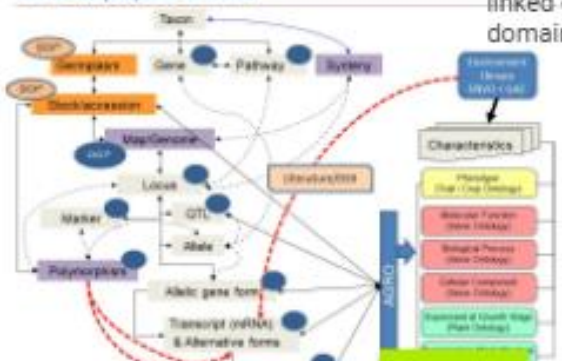
# INSDC record after injection of DOI

```
LOCUS      DQ884074                257 bp   mRNA    linear   EST 24-FEB-2011
DEFINITION DQ884074 Oryza sativa (indica cultivar-group) cv. IR64 cDNA-AFLP
           fragment Oryza sativa Indica Group cDNA clone 51_9b, mRNA sequence.
ACCESSION  DQ884074
VERSION    DQ884074.1
DBLINK     BioSample: SAMN00165158
KEYWORDS   EST.
SOURCE     Oryza sativa Indica Group (long-grained rice)
  ORGANISM Oryza sativa Indica Group
           Eukaryota; Viridiplantae; Streptophyta; Embryophyta; Tracheophyta;
           Spermatophyta; Magnoliopsida; Liliopsida; Poales; Poaceae; BOP
           clade; Oryzoideae; Oryzeae; Oryzinae; Oryza; Oryza sativa.
REFERENCE  1 (bases 1 to 257)
  AUTHORS  Ventelon-Debout,M., Tranchant-Dubreuil,C., Nguyen,T.-T.-H.,
           Bangratz,M., Sire,C., Delseny,M. and Brugidou,C.
  TITLE    Rice Yellow Mottle Virus stress responsive genes from susceptible
           and tolerant rice genotypes
  JOURNAL  BMC Plant Biol. 8, 26 (2008)
  PUBMED   18315879
COMMENT    Contact: Ventelon-Debout M
           UMR 5096
           Institut de Recherche pour le Developpement
           911 avenue d'Agropolis, BP54501, Montpellier, 34394, France.
FEATURES   Location/Qualifiers
  source   1..257
           /organism="Oryza sativa Indica Group"
           /mol_type="mRNA"
           /cultivar="IR64"
           /db_xref="taxon:39946"
           /db_xref="https://doi.org/10.18730/V918"
           /clone="51_9b"
           /clone_lib="SAMN00165158 Oryza sativa (indica
           cultivar-group) cv. IR64 cDNA-AFLP fragment"
```



GLIS DOI

What we propose to build



What key systems or types of systems we will need to connect with to be able to implement linked data management between HTP-enabled breeding and genebank operational domains--using DGOs/data interoperability as the 'connective tissue'?



Global Information System (GLIS) as a link to DOI System and onward link to Genesys/other systems

Germplasm that is publicly released gets DOI from GLIS

DGOs are generated for digital sequence information--and these are assigned DOIs.

Markers DGOs (in accomm different markers DGOs are kind of a Sequenc This can assigned







Platform for  
Big Data  
in Agriculture



**THANK YOU**

**[pankaj.jaiswal@oregonstate.edu](mailto:pankaj.jaiswal@oregonstate.edu)**