

Figure 4 : BIS architecture (database and interface layers)

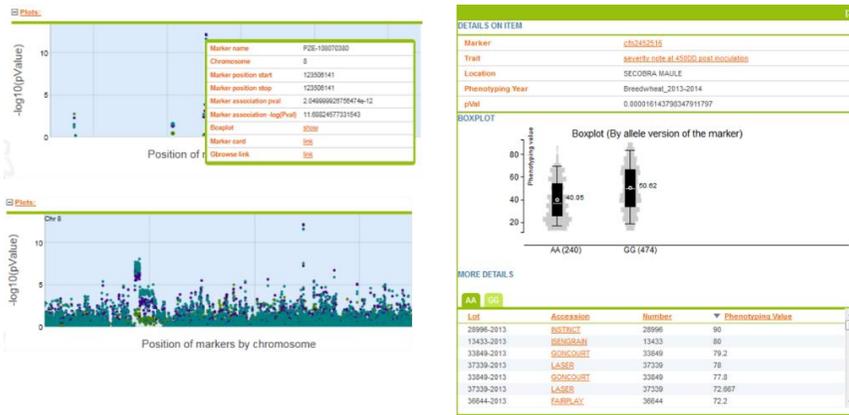


Figure 5 : GWAS graphics from BIS using the Google Web Toolkit technology to generate on-the-fly plots (Manhattan plot and Boxplot)

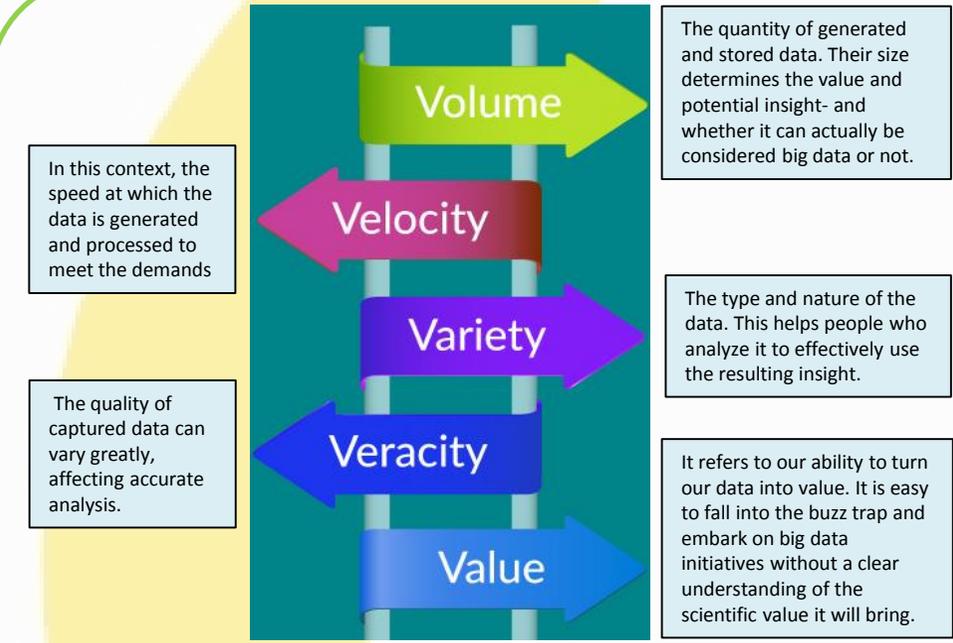


Figure 1: the 5Vs of Big Data

More information

www.breedwheat.fr

Coordinator : Dr. Jacques Le Gouis, UMR GDEC - jacques.le-gouis@inra.fr
Project manager : Emmanuelle Legendijk, INRA Transfert - emmanuelle.legendijk@paris.inra.fr
Communication : Bernard Bejar, Céréales Vallée - bernard.bejar@cereales-vallee.org

This project receives funding from the French Government managed by the Research National Agency (ANR) in the framework of the Investments for the Future (ANR-10-BTBR-03), France Agrimer and the French Fund to support Plant Breeding (FSOV).

The data generated by BreedWheat that are managed by bioinformaticians from WP5 (“Bioinformatics for gene discovery, data integration and dissemination”) can be considered as Big data because of their:

- Volume**: 3.7 billion genotypic data, 1 million phenotypic data (from 2011 to 2014), etc.
- Velocity**: new phenotypic data each year to generate and integrate in the Information System.
- Variety**: germplasm, SNP discovery, genetic maps, genotyping, phenotyping, ontology and genome wide association studies data.
- Veracity**: data are checked and curated by the bioinformaticians of the WP5.
- Value**: data analysis will lead to the development of new varieties.

The BreedWheat data

So the volume is quite big and the data are diverse, but more important, the data should be linked together, interoperable and displayed easily and quickly. The BreedWheat Information System (BIS) rely on the GnplS Information System (Steinbach *et al.*, GnplS: an information system to integrate genetic and genomic data from plants and fungi. Database (Oxford) 2013; 2013 bat058. doi: 10.1093/database/bat058). It allows the data integration and is chosen to host all BreedWheat data. BIS also manages the Intellectual Property of the data with some open data and some accessible only by the project partners before opening them to the community.

Breeding for economically and environmentally sustainable wheat varieties: an integrated approach from genomics to selection.



Data

BreedWheat Information System rely on GnplS with dedicated access rights following the consortium agreement.

- [Germplasm](#) (public access)
- [SNP discovery](#) : Major genes, Digital, BGA, BBSRC, Candidates, Infinium, IWGSC genes, IWGSC intergenic, Biogemma, ISBP.
- Axiom TaBW420K v2
 - Genotyping data
 - [TaBW420K_V2_BW_WP1_CsRe](#)
 - [TaBW420K_V2_BW_WP1_Discovery_plate](#)
 - [TaBW420K_V2_BW_WP1_ITMI_p1](#)
 - [TaBW420K_V2_BW_WP2_Core_Collection](#)
 - [TaBW420K_V2_BW_WP2_Elite_panel](#)
 - [TaBW420K_V2_BW_WP2_Extended_Elite_panel](#)
 - [TaBW420K_V2_BW_WP3_Genetic_Resources](#)
 - [TaBW420K_V2_BW_WP4_Genomic_Selection](#)
 - [TaBW420K_V2_BW_WP4_MAGIC](#)
 - Genetic map data
 - to [display in GnplS](#)
 - to [download](#)
 - Association data
 - in progress
- Phenotyping:
 - [INRA Small Grain Cereals Network Phenotypic Trials dataset](#) on 11 locations and 16 years for more than 1700 bread wheat genotypes (public access).
 - Dataset updated with [new 2015 data](#).
 - [BreedWheat](#)
 - 2011-2012 : 9 trials
 - 2012-2013 : 10 trials
 - 2013-2014 : 8 trials
- Ontology:
 - [Wheat INRA Phenotype Ontology \(WIPO\) v1](#) (public access)

Figure 2 : data integrated in the BreedWheat Information System (<https://wheat-urgi.versailles.inra.fr/Projects/BreedWheat>)

The data integration workflow

The data integration starts from data acquisition in the labs and the fields by the researchers and the technicians. Then the data are curated and formatted in dedicated Excel files (one per data type) including their meta-data (linked ontologies, data standards, etc.) with the help of the WP5 bioinformaticians.

The data managers use the Talend tool to extract, transform and load the data from the Excel files to the BIS Information System.

The development team adapts the Information System to the specific needs of the BreedWheat consortium and improves the databases and interfaces.

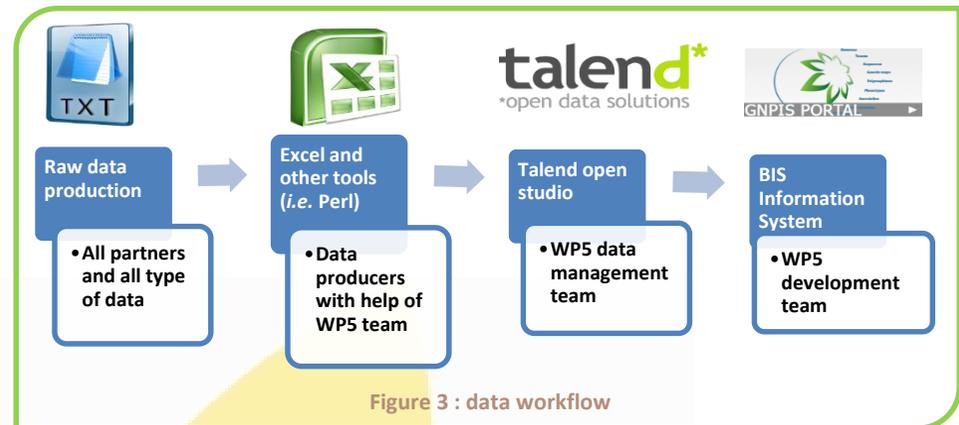


Figure 3 : data workflow

The new technologies

BIS is built on the following technologies:

- ❑ PostgreSQL v9.6 Relational database
- ❑ Elasticsearch v2.3 NoSQL database
- ❑ Google Web Toolkit (GWT) v2.7 framework

The database architecture is an hybrid system using PostgreSQL for the key data to insure the best data quality and consistency, and Elasticsearch for high volume data to allow their storage and fast data query.

BIS interface relies on the GWT framework to generate fast Javascript from Java code. In particular, INRA URGI uses dedicated libraries to display Genome Wide Association Studies (GWAS) plots (gwt.visualization) and synteny browser (gwt.graphics).