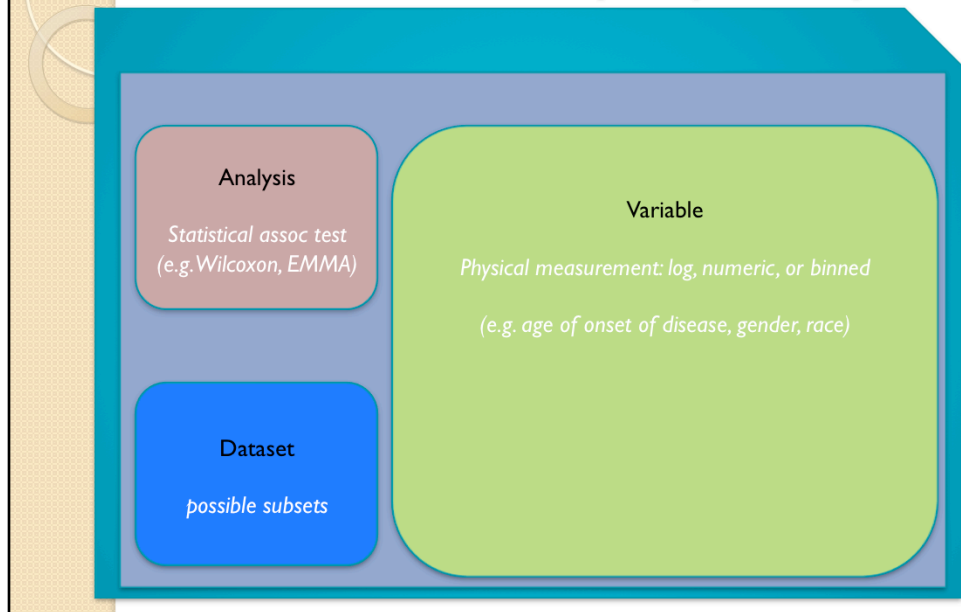# Gramene Retreat: A.t. Flowering Time GWAS

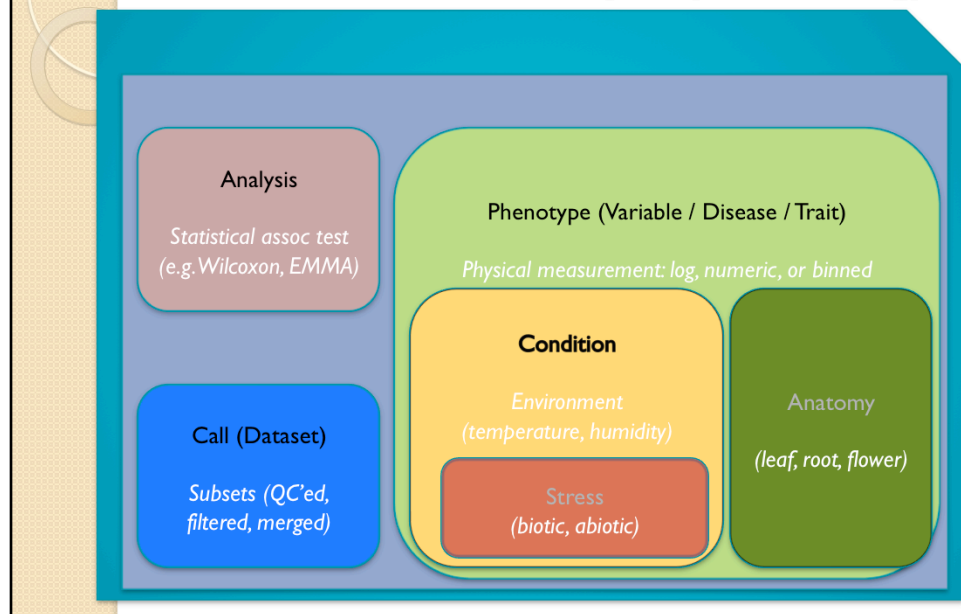

Chuah Aaron
17 June 2010

Gap#1: dbGaP is NCBI's pheno/genotypic equivalent of dbSNP – it stores Variations, Documents and Analyses (Associations), which may not be enough

Genome-wide association study of 107 phenotypes in *Arabidopsis thaliana inbred lines*

*Susanna Atwell, Yu S. Huang, Bjarni J. Vilhjálmsson, Glenda Willems, Matthew Horton, Yan Li, Dazhe Meng, Alexander Platt, Aaron M. Tarone, Tina T. Hu, Rong Jiang, N. Wayan Muliyati, Xu Zhang, Muhammad Ali Amer, Ivan Baxter, Benjamin Brachi, Joanne Chory, Caroline Dean, Marilyne Debieu, Juliette de Meaux, Joseph R. Ecker, Nathalie Faure, Joel M. Kniskern, Jonathan D. G. Jones, Todd Michael et al.*

Gap #2: You've seen the verbosity of the Disease/Trait (Phenotype) naming yesterday
http://uswest.ensembl.org/Homo_sapiens/Variation/Phenotype?
source=dbSNP;v=rs420259

Gap #3: Nordborg's data contains much information, but without an appropriate visualization, it's hard to interpret

# DAS tracks

Implemented DAS Tracks:

❖ **Histogram**
- ❖ p-value (-log10 transformed)
- ❖ (phenotype) risk-allele frequency (0 to 0.5)
- ❖ rank (#1 to #1000), plotted in reverse (highest bar=highest rank)

**Grouping & Coloring:**
- ❖ Phenotypes of the same class (Flowering time, Bolt length 5cm, etc)
  - ❖ Different Ions/Stresses (Arsenic, Cu, Co, Zi, etc)
  - ❖ Temperature conditions (10°C, 16°C, 22°C)
  - ❖ Analysis method (EMMA, KW, etc)

http://uswest.ensembl.org/Homo_sapiens/Variation/Phenotype?
source=dbSNP;v=rs420259

## Phenotype Categories:

- ❖Flowering Time (between germination & emergence of first flowers)
- ❖Flowering Gene (FRI & FLC RNA gene expression levels)
- ❖Leaf number, curl, roll & serration
- ❖Development (time for bolt to reach height of 5cm)
- ❖Germination (time to emergence of cotyledons)
- ❖Seed Storage (primary & secondary dormancy)
- ❖Width (average diameter of 4 plants 8 weeks after germ)
- ❖Collapse (of leaves after inoculation with 0.1ml of bacteria)
- ❖Other (other inoculations / ICP-MS ion concentrations)
  - ❖Arsenic, Boron, Calcium, Cadmium, Cobalt, Copper, Iron, Potassium, Lithium, Magnesium, Manganese, Molybdenum, Nickel, Phosphorus, Sodium, Sulfur, Selenium, Zinc

http://uswest.ensembl.org/Homo_sapiens/Variation/Phenotype?source=dbSNP;v=rs420259

Gene: Frigida

Flowering Locus C

Note low rank of red line near 3.19M despite p-value close to the rest (and higher maf than most)

Primarily Red (FT), Pink (FT-gxp) & Blue (Development)

# Phenotype description

255 chars isn't enough to capture Nordborg's "growth_condition" & "phenotype_scoring" fields

- e.g. phenotype "avrPphB"
  - 20degC, 12 hrs daylight
  - Following inoculation of two leaves per plant with 0.1 ml of 10 -8 cfu/ml bacteria in 10 mM MgSO4 buffer using a blunt-tipped syringe, leaf collapse was scored at 20 hrs and again at 24 hrs after inoculation. A positive score at either time point was deemed a hypersensitive response
  - Proposed shortening
    - "hypersensitivity: 12h-day 20C, P.syringae AvrPphB"

# Other examples

- ❖ "SDV"
  - ❖ 18 degC, 8 hrs daylight, vernalized (5 wks, 4 degC)
  - ❖ Number of days following stratification to opening of first flower. The experiment was stopped at 200 d, and accessions that had not flowered at that point were assigned a value of 200
  - ❖ Proposed shortening
    - ❖ "flowering time: 8h-day 18C, vernalized 5w 4C"
- ❖ "At1 CFU2 "
  - ❖ 20 degC, 12 hrs daylight
  - ❖ In planta bacterial growth (number of CFU / leaf area) of the P. viridiflava strain was individually measured as described in Goss and Bergelson 2006
  - ❖ Proposed shortening
    - ❖ "leaf CFU: 12h-day 20C, P.viridiflava At1"

# EO:Environment Ontology mapping

**Environmental condition**
- Duration of conditioning:
  - 8w, 5w, unspecified
- Day length:
  - 8h (short), 16h (long), 10h, unspecified, 0h (dark)
- Temperature:
  - 10C, 16C, 22C, 20-22C, 23C
- Humidity:
  - 50% hu, 70% hu
- Location:
  - Field, Greenhouse
  - Dry Storage

**[+ Stress condition]**
- Vernalization (cold regimen):
  - duration, day length, temperature, humidity
- Biotic
  - viral, bacterial strains

http://uswest.ensembl.org/Homo_sapiens/Variation/Phenotype?source=dbSNP;v=rs420259

# TO: Trait Ontology mapping

**Growth trait**
- Flowering Time
  - Flowering Duration (last-first flower)
  - Bolt time (to 5cm)
- Senescence (life cycle)
  - Reproduction Time (senes-first flower)
  - Maturation Time (senes-flower senes)
- Germination Time

**Stress trait**
- Biotic
  - (Symptom)
  - colony forming units (CFU)
  - Aphid number
- Abiotic
  - vernalization

**Anatomy/Morphology trait**
- Leaf
  - Number
  - Shape
    - Serration, rolled, rosette erect
  - (Symptom)
    - Chlorosis, anthocyanin, lesioning (necrosis)
- Stem
  - Diameter
- Root
  - Trichome density
- Seed
  - Hypocotyl length

**Biochemical trait**
- Ion Concentration
  - Li, B, Na, Mg, P, S, K, Ca, Mn, Fe, Co, Ni, Cu, Zn, As, Se, Mo, Cd

http://uswest.ensembl.org/Homo_sapiens/Variation/Phenotype?
source=dbSNP;v=rs420259

# Controlled Vocabulary Phenotypes

avrPphB

❖ leaf collapse: 12h-day 20C, P.syringae AvrPphB
  ❖ Hypersensitivity response <GO:0002524>
  ❖ Bacterial disease resistance <TO:0000315 P.syringae AvrPphB>
  ❖ warm/hot temperature regimen <EO:0007173 20C>
  ❖ intermittent light regimen <EO:0007128 12h>

SDV

❖ flowering time: 8h-day 18C, vernalized 5w 4C

At1 CFU2

❖ leaf CFU: 12h-day 20C, P.viridiflava At1

To send to Josh to work out the first pass...
*AvrPphB* is an avirulence (Avr) protein from the plant pathogen Pseudomonas syringae that can trigger a disease-resistance response

# Data Query Entry Points / Use cases

- ❖ SNP quality/confidence
  - ❖ Subsets of SNPs (variation set functionality in Ensembl r58)
    - ❖ 470k (submitted), 71k (genotyped), 9k (high-confidence) for grape
- ❖ Entry points for variation Mart / new form to get sets of SNPs
  - ❖ by position (genomic range)
  - ❖ by confidence (SNP set)
  - ❖ by sample accession (ecotype)
    - ❖ Difference between ref strand and specific accession
    - ❖ Difference between 2 accessions
      - ❖ sites within a specified region where Mo17 and CML vary
  - ❖ by function (consequence)
  - ❖ by annotation
    - ❖ Within 10kb upstream of a given gene list
    - ❖ Associated with flowering time phenotypes
      - ❖ With a p-value of <1E-4
  - ❖ by Linkage Disequilibrium (LD) threshold
  - ❖ by allele Frequency

Gap #4: more entry points needed

# Arabidopsis 2010 Viz roadmap

- ❖ Google Web Toolkit (GWT) 2.1 based view (end-June)
  - ❖ Plugs in into any populated Gramene/Ensembl Var DB
  - ❖ to benchmark alongside GWT Data Presentation Widgets
  - ❖ Variation_annotation selection/subclassing/grouping
    - ❖ By study
    - ❖ By phenotype (i.e. flowering time)
    - ❖ By condition (i.e. long day = 16h-day, vernalization)
    - ❖ By sample (ecotype)
      - ❖ i.e. the variations of that sample that differ from the reference
  - ❖ Ability to store whether the association value is a
    - ❖ P-value, Score, Rank
    - ❖ Log-transformation used (if any)
    - ❖ And should the actual sample-phenotype measurement be stored

Pulling from GDPDM?