# Genome-wide analysis of a wheat transcription factor family:

# The power of bioinformatics resources

**27th May 2020**

Susanne Schilling & Rainer Melzer
School of Biology and Environmental Science
University College Dublin

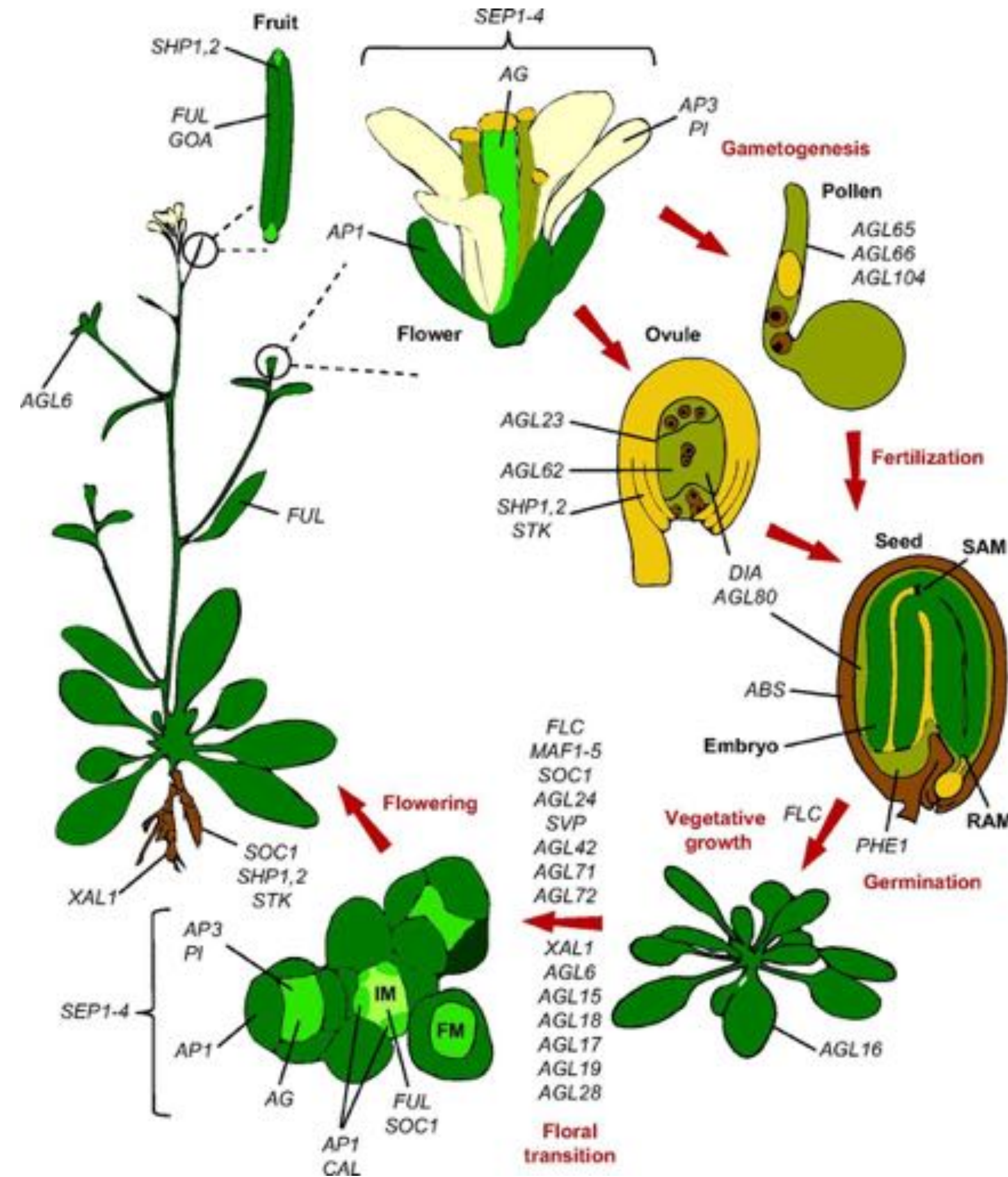@UCDflowerpower          https://ucdflowerpower.org/

# What drives domestication?

Domestication genes: genes with allelic versions that contributed to cultivating plants for human needs.

3 main groups of genes/proteins:

1. enzymes or structural proteins    "superheros"

2. numerous genes collectively    "minions"
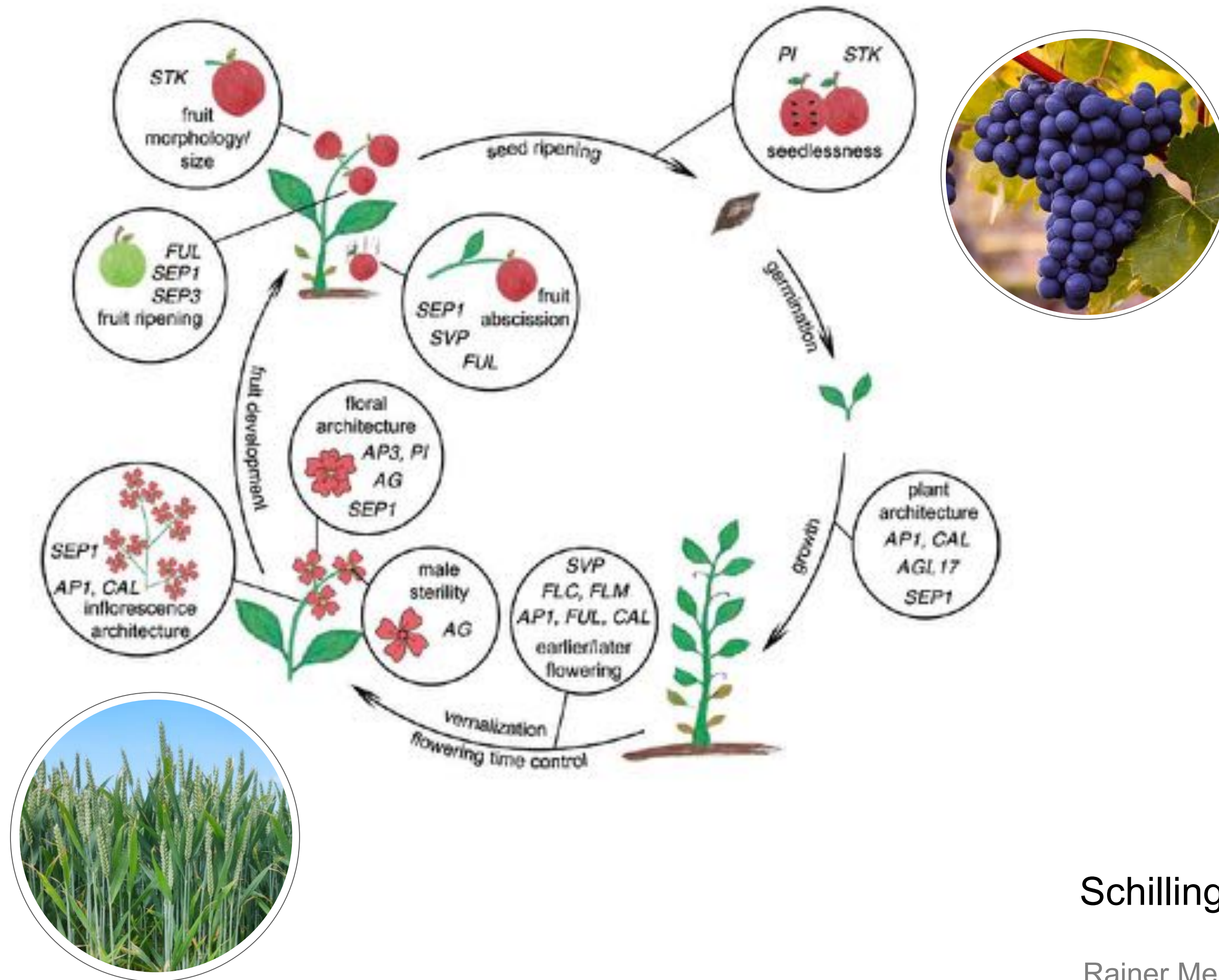
3. transcription factors    "masterminds"

Martínez-Ainsworth & Tenaillon, 2016
Schilling et al., 2018

Rainer Melzer & Susanne Schilling, UCD, Ireland

# MADS-box genes are key regulators of plant development



Smaczniak et al., 2012

Rainer Melzer & Susanne Schilling, UCD, Ireland

# MADS-box genes and plant domestication



Schilling et al., 2018

Rainer Melzer & Susanne Schilling, UCD, Ireland

# Genome-wide characterization of MADS-box genes in wheat

Rainer Melzer & Susanne Schilling, UCD, Ireland

# Why wheat is not a model plant

**diploid (*Arabidopsis* and rice)**

2n

homologs

1 gene
different alleles

**hexaploid (wheat)**

homologs

6n

homoeologs

Marcussen et al., 2014

Rainer Melzer & Susanne Schilling, UCD, Ireland

# Why wheat is not a model plant



Genome size in million base pairs

# Why wheat is not a model plant

Genome size in million base pairs

Rainer Melzer & Susanne Schilling, UCD, Ireland

# With the help of a high quality wheat genome



RESEARCH

**RESEARCH ARTICLE**

WHEAT GENOME

## Shifting the limits in wheat research and breeding using a fully annotated reference genome

International Wheat Genome Sequencing Consortium (IWGSC)*

Rainer Melzer & Susanne Schilling, UCD, Ireland

# MADS-box genes are highly conserved in wheat



Schilling et al., 2020

Rainer Melzer & Susanne Schilling, UCD, Ireland

# MADS-box genes are highly conserved in wheat



http://bar.utoronto.ca/eplant_wheat/

Schilling et al. 2020

Rainer Melzer & Susanne Schilling, UCD, Ireland

# MADS-box genes are highly conserved in wheat

Schilling et al. 2020



Rainer Melzer & Susanne Schilling, UCD, Ireland

# Putative neofunctionalization of MADS-box genes?



Schilling et al. 2020

Rainer Melzer & Susanne Schilling, UCD, Ireland

# Some MADS-box gene subfamilies underwent expansion in wheat



Schilling et al. 2020

Rainer Melzer & Susanne Schilling, UCD, Ireland

# Some MADS-box gene clades have undergone expansion in wheat – how?



Recombination rate

Gene density

Choulet et al., 2014

Rainer Melzer & Susanne Schilling, UCD, Ireland

# Chromosomal distribution of wheat MADS-box genes



Schilling et al. 2020

Rainer Melzer & Susanne Schilling, UCD, Ireland

# Subtelomeric segments as hot spots for MADS-box gene evolution?



Schilling et al. 2020

Rainer Melzer & Susanne Schilling, UCD, Ireland

# Subtelomeric segments as hot spots for MADS-box gene evolution?



$y = 1.7176x + 4.4631$
$R^2 = 0.61353$

Putative neo-functionalization

Conserved core functions

Schilling et al. 2020

Rainer Melzer & Susanne Schilling, UCD, Ireland

# Summary I

- MADS-box genes are key players in wheat development

- Conserved sequence and expression pattern

- MADS-box genes might have contributed to the success of wheat by neo- and subfunctionalization

- Good candidates for crop improvement

Rainer Melzer & Susanne Schilling, UCD, Ireland

# Genome-wide analysis of a wheat transcription factor family:

## The power of bioinformatics resources

# Part II: a look under the hood

bioinformatics

???

wheat
resources

# From Genes to Phylogenies

gene mining and identification

↓

filtering and sorting

↓

phylogeny



Schilling et al. 2020

Rainer Melzer & Susanne Schilling, UCD, Ireland

# From Genes to Phylogenies

**all IWGSC genes**

select by PFAM domain
(PF00319; PF01486)

**all genes MADS/K**

gene mining and identification

↓

filtering and sorting

↓

phylogeny

## Index of /download/iwgsc/IWGSC_RefSeq_Annotations/v1.1

| [ICO] | Name | Last modified | Size | Description |
|-------|------|---------------|------|-------------|
| [DIR] | Parent Directory | | - | |
| [ ] | iwgsc_refseqv1.1_README.pdf | 20-Aug-2019 11:31 | 14K | |
| [ ] | iwgsc_refseqv1.1_genes_2017July06.zip | 05-Dec-2017 16:36 | 283M | |
| [TXT] | iwgsc_refseqv1.1_genes_2017July06.zip.md5.txt | 05-Dec-2017 16:36 | 70 | |
| [ ] | iwgsc_refseqv1.1_manually_curated_gene_families.zip | 12-Oct-2018 18:40 | 282K | |
| [TXT] | iwgsc_refseqv1.1_manually_curated_gene_families.zip.md5.txt | 12-Oct-2018 18:40 | 84 | |
| [ ] | iwgsc_refseqv1.1_rnaseq_mapping_2017July20.zip | 05-Dec-2017 16:30 | 1.6G | |
| [TXT] | iwgsc_refseqv1.1_rnaseq_mapping_2017July20.zip.md5.txt | 05-Dec-2017 16:35 | 79 | |

https://urgi.versailles.inra.fr/download/iwgsc/IWGSC_RefSeq_Annotations/v1.1/

Rainer Melzer & Susanne Schilling, UCD, Ireland

# From Genes to Phylogenies

all IWGSC genes

select by PFAM domain
(PF00319; PF01486)

gene mining and identification

all genes MADS/K

filtering and sorting

phylogeny

# IWGSC Wheat mine

Rainer Melzer & Susanne Schilling, UCD, Ireland

# IWGSC Wheat mine

Rainer Melzer & Susanne Schilling, UCD, Ireland

# IWGSC Wheat mine

# Data Mining - IWGSC Wheat mine

# Ensembl Plants – Data Mining



plants.ensembl.org/Triticum_aestivum/

Rainer Melzer & Susanne Schilling, UCD, Ireland

# Ensembl Plants – Data Mining



plants.ensembl.org/Triticum_aestivum/

Rainer Melzer & Susanne Schilling, UCD, Ireland

# Ensembl Plants – Data Mining



- FASTA

# Ensembl Plants – Data Mining



plants.ensembl.org/Triticum_aestivum/

Rainer Melzer & Susanne Schilling, UCD, Ireland

# From Genes to Phylogenies

all IWGSC genes

select by PFAM domain
(PF00319; PF01486)

gene mining and identification

all genes MADS/K

BLAST

filtering and sorting

phylogenetic

MIKC-type MADS

MAFFT version 7
Multiple alignment program for amino acid o

https://mafft.cbrc.jp/alignment/server/

phylogeny

Rainer Melzer & Susanne Schilling, UCD, Ireland

# MAFFT for building alignments

Rainer Melzer & Susanne Schilling, UCD, Ireland

# From Genes to Phylogenies

**all IWGSC genes**

select by PFAM domain
(PF00319; PF01486)

gene mining and identification

**all genes MADS/K**

BLAST

filtering and sorting

phylogenetic

**MIKC-type MADS**

phylogeny

e! EnsemblPlants

WheatMine

URGI

e! EnsemblPlants

MAFFT version 7
Multiple alignment program for amino acid o

https://mafft.cbrc.jp/alignment/server/

IQ-TREE

http://www.iqtree.org/

UCD DUBLIN

Rainer Melzer & Susanne Schilling, UCD, Ireland

# From Genes to Phylogenies



**all IWGSC genes**

select by PFAM domain
(PF00319; PF01486)

**WheatPlants**

**WheatMine**

gene mining and identification

**all genes MADS/K**

URGI

**BLAST**

EnsemblPlants

filtering and sorting

phylogenetic

MAFFT

IQ-TREE

**MIKC-type MADS**

phylogeny

curating

**MIKC-type MADS**

Rainer Melzer & Susanne Schilling, UCD, Ireland

# From Genes to Phylogenies



all IWGSC genes

gene mining and identification

select by PFAM domain
(PF00319; PF01486)

all genes MADS/K

BLAST

filtering and sorting

phylogenetic

MIKC-type MADS

phylogeny

curating

CDD

FGENESH+

MIKC-type MADS

Rainer Melzer & Susanne Schilling, UCD, Ireland

# From Genes to Phylogenies



Rainer Melzer & Susanne Schilling, UCD, Ireland

# From Genes to Phylogenies



all IWGSC genes

select by PFAM domain
(PF00319; PF01486)

gene mining and identification

all genes MADS/K

BLAST

filtering and sorting

phylogenetic

MIKC-type MADS

phylogeny

FGENESH+

curating

CDD

MIKC-type MADS

Rainer Melzer & Susanne Schilling, UCD, Ireland

# Phylogeny tools

MAFFT version 7
Multiple alignment program for amino acid or nucleotide sequences

https://mafft.cbrc.jp/alignment/server/

IQ-TREE

http://www.iqtree.org/

geneious

https://www.geneious.com/

FigTree v1.4.4

https://github.com/rambaut/figtree/releases



**MIKC-type MADS**

wheat   201

rice   43

Arabidopsis   45

Schilling et al. 2020

Rainer Melzer & Susanne Schilling, UCD, Ireland

# Counting, comparing and stats



Schilling et al. 2020

Rainer Melzer & Susanne Schilling, UCD, Ireland

# Circos plots to visualize whole genomes

Rainer Melzer & Susanne Schilling, UCD, Ireland

# Circos plot with Shiny Circos



shinyCircos: an R/Shiny application for interactive creation of Circos plot

https://github.com/venyao/shinyCircos

http://150.109.59.144:3838/shinyCircos/

Rainer Melzer & Susanne Schilling, UCD, Ireland

# Circos plot with Shiny Circos



https://github.com/venyao/shinyCircos

http://150.109.59.144:3838/shinyCircos/

Rainer Melzer & Susanne Schilling, UCD, Ireland

# Expression analysis



Data mining

Visualisation

Schilling et al. 2020

Rainer Melzer & Susanne Schilling, UCD, Ireland

# Expression data

Rainer Melzer & Susanne Schilling, UCD, Ireland

# Genevestigator



https://genevestigator.com/

Rainer Melzer & Susanne Schilling, UCD, Ireland

# Genevestigator



Rainer Melzer & Susanne Schilling, UCD, Ireland

https://genevestigator.com/

# Wheat-expression.com



- TPM
- Counts
- PNG

http://www.wheat-expression.com/

Rainer Melzer & Susanne Schilling, UCD, Ireland

# Heatmap tools

Rainer Melzer & Susanne Schilling, UCD, Ireland

# More wheat resources…



https://urgi.versailles.inra.fr/jbrowseiwgsc/gmod_jbrowse/



https://urgi.versailles.inra.fr/synteny/synteny/



*Aegilops tauschii*
*Tritcum urartu*

Synteny
EMS mutants

plants.ensembl.org/

*Aegilops speltoides*

http://wheat-urgi.versailles.inra.fr/

# Other resources…

https://www.ncbi.nlm.nih.gov/

http://hmmer.org/

https://www.youtube.com/

https://rstudio.com/

**Galaxy**

https://usegalaxy.org/

https://www.biostars.org/

https://www.edx.org/

MAFFT version 7
Multiple alignment program for amino acid or nucleotide sequences

https://mafft.cbrc.jp/alignment/server/

START: Shiny Transcriptome Analysis Resource Tool

https://kcvi.shinyapps.io/START/

EMBL-EBI Hinxton

https://www.ebi.ac.uk/training/online/

IQ-TREE

FigTree v1.4.4

http://www.iqtree.org/

https://github.com/rambaut/figtree/releases

Rainer Melzer & Susanne Schilling, UCD, Ireland

Thank you for listening

Questions?

SiRui Pan

Lars Jermiin

Alice Kennedy

susanne.schilling@ucd.ie

rainer.melzer@ucd.ie

Susanne Schilling & Rainer Melzer
School of Biology and Environmental Science
University College Dublin