

# **Introduction to Machine Learning for Genomics**

Research examples

Director and Professor

Department of Evolutionary Biology and Environmental  
Studies

University of Zurich

Kentaro Shimizu

[kentaro.shimizu@uzh.ch](mailto:kentaro.shimizu@uzh.ch)

# Two published examples of genomic researches using machine learning

- Genome-wide association studies and genomic prediction of biodiversity effects
  - LASSO regression both for choosing interpretable variables and for prediction
- Hardware and software PlantServation for time-course image analysis of polyploid species in field
  - Deep learning, random forest

# Examples of NGS analysis

- *de novo* genome assembly
- genome-wide polymorphism patterns, selection scan, and genome-wide association studies (GWAS)
- RNA-seq (expression analysis, transcriptome)
  - cDNA assembly
- metagenomics
- small RNA
- ChIP-Seq

# NGS analysis: assembly and mapping

## Assembly

DNA: genome assembly

Keyword: contig, coverage

## Mapping to a reference genome/transcriptome

DNA: polymorphism analysis and genome-wide association studies

    DNA methylation, ChIP-seq

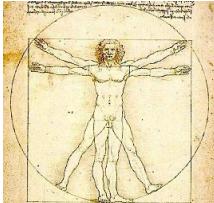
RNA: RNA-seq

28 Oct 2010: 179 human genomes were reported. How are they meaningful?



# Why do genome-wide polymorphisms matter?

## Projects for medical, agricultural, ecological & evolutionary biology

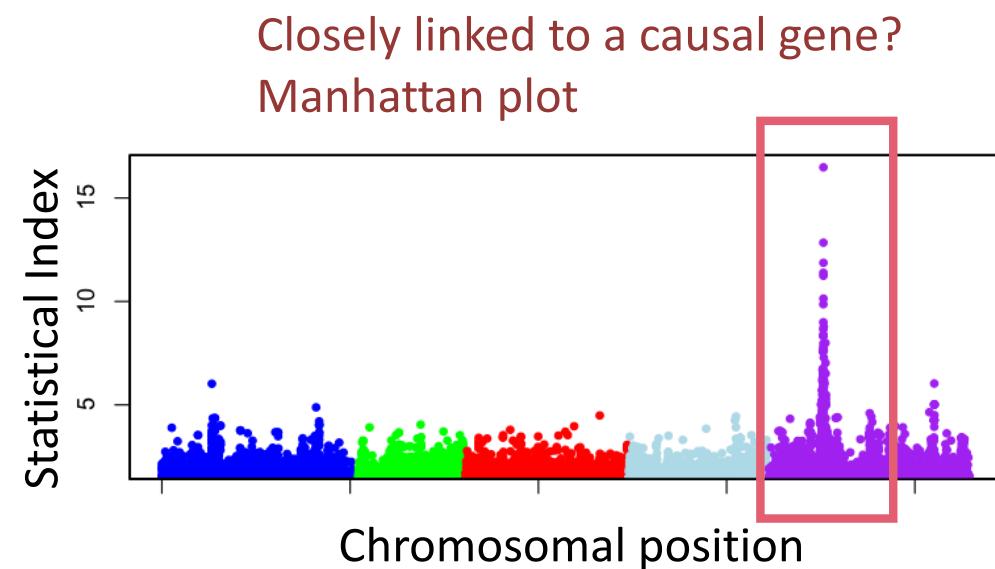
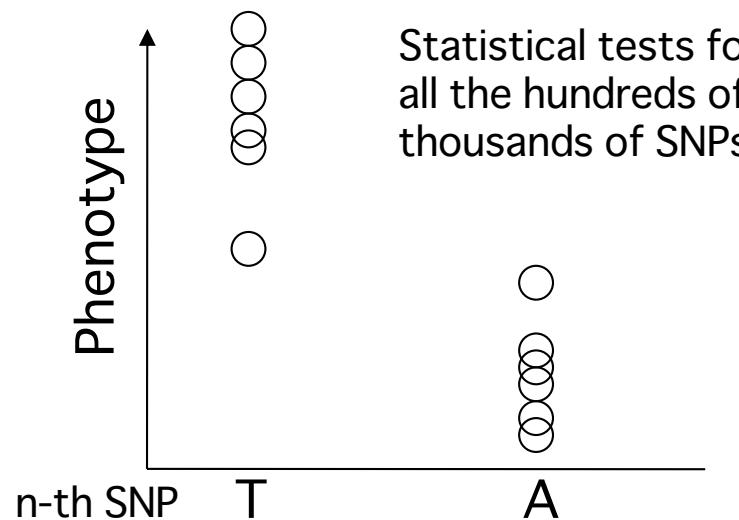
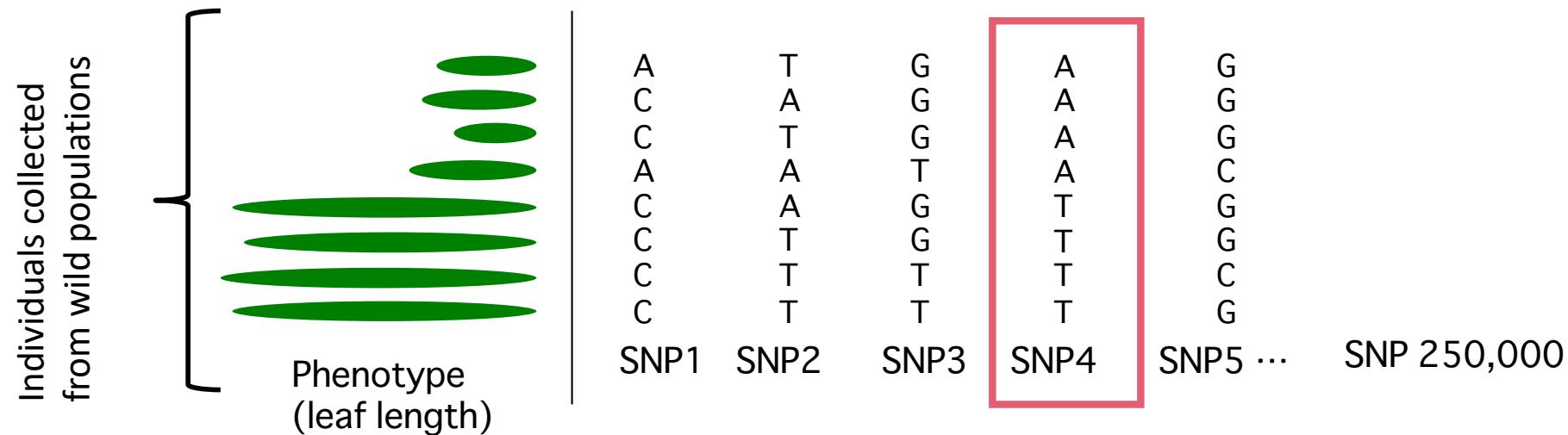


	<i>Homo sapiens</i>	<i>Drosophila melanogaster</i>	<i>Arabidopsis thaliana</i>	<i>Triticum aestivum</i> (bread wheat)
Size genome	3 Gb	180 Mb	130 MB	16 Gb, allohexaploid (6x)
SNP (single nucleotide polymorphism)	HapMap <i>Nature</i> 2005;2007		96 accessions, 876 loci, Sanger Nordborg et al. <i>PLoS Biol</i> 2005 250 k SNPs, 191 accessions Kim et al. <i>Nat Genet</i> 2007, Atwell et al. <i>Nature</i> 2010	
re-sequencing or de novo	C. Venter J. Watson 1000 genomes	192 genomes DGRP <i>Nature</i> 2012	20 genomes, Clark et al., <i>Science</i> 2007 1001 genomes (Gan et al. <i>Nature</i> 2011, <i>Cell</i> 2016)	10+ genomes (de novo) (Walkowiak et al. <i>Nature</i> 2020; Shimizu et al. <i>PCP</i> 2021)
related species	Chimp Neanderthal	Drosophila 12 Genomes	<i>A. lyrata</i> (Hu et al. <i>Nature Genetics</i> 2011), <i>Cardamine hirsuta</i> (Gan et al. <i>Nature Plants</i> 2016), <i>A. halleri</i> (Akama et al. <i>NAR</i> 2014; Briskine et al. 2016) Phylogenomics consortium (Novikova et al. <i>Nature Genet</i> 2016)	4x wheat (Avni et al. <i>Science</i> 2017) and many others

# Genome-Wide Association Study

Developed for human disease studies

Looks for the genes responsible for the phenotypic variations based on simple associations between phenotypes and genotypes (e.g. SNPs)



# GWAS important for medical studies

Weedon et al. Nature Genet 39:1245, 2007

Human height is a classic, highly heritable quantitative trait. To begin to identify genetic variants influencing height, we examined genome-wide association data from 4,921 individuals. Common variants in the *HMGA2* oncogene, exemplified by rs1042725, were associated with height ( $P = 4 \cdot 10^{-8}$ ). *HMGA2* is also a strong biological candidate for height, as rare, severe mutations in this gene alter body size in mice and humans, so we tested rs1042725 in additional samples. We confirmed the association in 19,064 adults from four further studies ( $P = 3 \cdot 10^{-11}$ , overall  $P = 4 \cdot 10^{-16}$ , including the genome-wide association data). We also observed the association in children ( $P = 1 \cdot 10^{-6}$ ,  $N = 6,827$ ) and a tall/short case-control study ( $P = 4 \cdot 10^{-6}$ ,  $N = 3,207$ ). We estimate that rs1042725 explains 0.3% of population variation in height (0.4 cm increased adult height per C allele). There are few examples of common genetic variants reproducibly associated with human quantitative traits; these results represent, to our knowledge, the first consistently replicated association with adult and childhood height.

**HMGA2 gene:  
mutation in mouse results in pygmy:**

Vol 447 | 28 June 2007 | doi:10.1038/nature05887

nature

## ARTICLES

### Genome-wide association study identifies novel breast cancer susceptibility loci

Douglas F. Easton<sup>1</sup>, Karen A. Pooley<sup>2</sup>, Alison M. Dunning<sup>3</sup>, Paul D. P. Pharoah<sup>2</sup>, Deborah Thompson<sup>1</sup>, Dennis G. Ballinger<sup>2</sup>, Jeffery P. Struwing<sup>4</sup>, Jonathan Morrison<sup>2</sup>, Helen Field<sup>2</sup>, Robert Luben<sup>2</sup>, Nicholas Wareham<sup>5</sup>, Shahana Ahmed<sup>2</sup>, Catherine S. Healey<sup>2</sup>, Richard Bowman<sup>6</sup>, the SEARCH collaborators<sup>2\*</sup>, Kerstin B. Meyer<sup>7</sup>, Christopher A. Hallman<sup>8</sup>, Laurence K. Kolonel<sup>9</sup>, Brian E. Henderson<sup>10</sup>, Loïc Le Marchand<sup>10</sup>, Paul Brennan<sup>10</sup>, Suleeporn Sangrajrang<sup>11</sup>, Valerie Gaborieau<sup>12</sup>, Fabrice Odefrey<sup>13</sup>, Chen-Yang Shen<sup>14</sup>, Pei-Ei Wu<sup>15</sup>, Hui-Chun Wang<sup>15</sup>, Diana Eccles<sup>15</sup>, D. Gareth Evans<sup>14</sup>, Julian Peto<sup>15</sup>, Olivia Fletcher<sup>16</sup>, Nicola Johnson<sup>16</sup>, Sheila Seal<sup>17</sup>, Michael R. Stratton<sup>17,18</sup>, Nazneen Rahman<sup>17</sup>, Georgia Chenevix-Trench<sup>19</sup>, Stig E. Bojesen<sup>20</sup>, Berge G. Nordestgaard<sup>20</sup>, Christen K. Axelsson<sup>21</sup>, Montserrat Garcia-Closas<sup>22</sup>, Louise Brinton<sup>22</sup>, Stephen Chanock<sup>23</sup>, Jolanta Lissowska<sup>24</sup>, Beata Peplonska<sup>25</sup>, Heli Nevanlinna<sup>26</sup>, Rainer Fagerholm<sup>26</sup>, Hannaleena Eraola<sup>26,27</sup>, Daehae Kang<sup>28</sup>, Keun-Young Yoo<sup>28,29</sup>, Dong-Young Noh<sup>28</sup>, Sei-Hyun Ahn<sup>30</sup>, David J. Hunter<sup>31,32</sup>, Susan E. Hankinson<sup>32</sup>, David G. Cox<sup>31</sup>, Per Hall<sup>33</sup>, Sara Wedren<sup>33</sup>, Jianjun Liu<sup>34</sup>, Yen-Ling Low<sup>34</sup>, Natalia Bogdanova<sup>35,36</sup>, Peter Schürmann<sup>36</sup>, Thilo Dörk<sup>36</sup>, Rob A. E. M. Tollenaar<sup>37</sup>, Catharine E. Jacob<sup>38</sup>, Peter Devilee<sup>38</sup>, Jan G. M. Klijn<sup>39</sup>, Alice J. Sigurdson<sup>40</sup>, Michael M. Doody<sup>41</sup>, Bruce H. Alexander<sup>42</sup>, Jinghui Zhang<sup>43</sup>, Angela Cox<sup>43</sup>, Ian W. Brock<sup>43</sup>, Gordon MacPherson<sup>43</sup>, Malcolm W. R. Reed<sup>43</sup>, Fergus J. Couch<sup>44</sup>, Ellen L. Goode<sup>45</sup>, Janet E. Olson<sup>45</sup>, Hanne Meijers-Heijboer<sup>44,47</sup>, Ans van den Ouwendijk<sup>44</sup>, André Uitterlinden<sup>48</sup>, Fernando Rivadeneira<sup>49</sup>, Roger L. Milne<sup>49</sup>, Gloria Ribas<sup>49</sup>, Anna Gonzalez-Neira<sup>49</sup>, Javier Benitez<sup>49</sup>, John L. Hopper<sup>50</sup>, Margaret McCredie<sup>51</sup>, Melissa Southey<sup>50</sup>, Graham G. Giles<sup>52</sup>, Chris Schoen<sup>53</sup>, Christina Justenhoven<sup>54</sup>, Hiltrud Braucht<sup>54</sup>, Ute Hamann<sup>55</sup>, Yon-Dschun Ko<sup>55</sup>, Amanda B. Spurdle<sup>59</sup>, Jonathan Beesley<sup>59</sup>, Xiaoqing Chen<sup>59</sup>, kConFab<sup>57,58</sup>, AOCS Management Group<sup>59</sup>, Arto Mannemaa<sup>58,59</sup>, Veli-Matti Kosma<sup>58,59</sup>, Vesa Kataja<sup>58,60</sup>, Jaana Hartikainen<sup>58,59</sup>, Nicholas E. Day<sup>2</sup>, David R. Cox<sup>2</sup> & Bruce A. J. Ponder<sup>2,7</sup>



Figure 9-1 Evolutionary Analysis, 4/e

© 2007 Pearson Prentice Hall, Inc.

Wood et al. Nature Genet 46:1173, 2014

### Defining the role of common variation in the genomic and biological architecture of adult human height

Using genome-wide data from 253,288 individuals, we identified 697 variants at genome-wide significance that together explained one-fifth of the heritability for adult height. By testing different numbers of variants in independent studies, we show that the most strongly associated ~2,000, ~3,700 and ~9,500 SNPs explained ~21%, ~24% and ~29% of phenotypic variance. Furthermore, all common variants together captured 60% of heritability. The 697 variants clustered in 423 loci were enriched for genes, pathways and tissue types known to be involved in growth and together implicated genes and pathways not highlighted in earlier efforts, such as signaling by fibroblast growth factors, WNT/β-catenin and chondroitin sulfate-related genes. We identified several genes and pathways not previously connected with human skeletal growth, including mTOR, osteoglycin and binding of hyaluronic acid. Our results indicate a genetic architecture for human height that is characterized by a very large but finite number (thousands) of causal variants.

### Genome-Wide Association Analysis Identifies Loci for Type 2 Diabetes and Triglyceride Levels

Diabetes Genetics Initiative of Broad Institute of Harvard and MIT, Lund University, and Novartis Institutes for BioMedical Research†

Science 316, 1331 (2007)

# Genome-wide association studies and genomic prediction of biodiversity effect using machine learning

Sato, Shimizu-Inatsugi, Takeda, Schmid, Nagano, Shimizu

*Nature Communications* 15: 8467, 2024

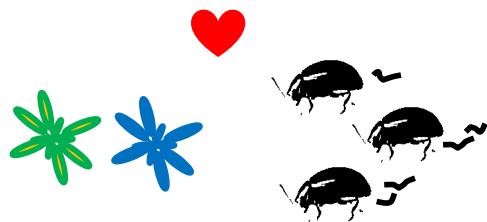
Similar to genomic prediction:  
Polygenic Risk Score (PRS) in medicine

The screenshot shows the University of Zurich (UZH) website. At the top, there is a navigation bar with links for "All News", "For Media", "UZH Magazin", and "Social Media". Below the navigation bar, a breadcrumb trail shows the path: "Home / All articles / Schaedlinge". A timestamp "07.10.2024 | Plant protection" is also present. The main headline reads "Reducing Herbivore Damage Using Biodiversity Instead of Insecticide". A brief summary follows: "Pesticides aren't always necessary: researchers at the University of Zurich have conducted a comprehensive field study showing that damage from herbivores can be reduced by using biodiversity within a plant species. Different plant genotypes can cooperate to help fend off herbivorous insects." At the bottom of the page, there are links to "Mathematics and Natural Sciences", "Sustainability", "Research", and "Media Releases".

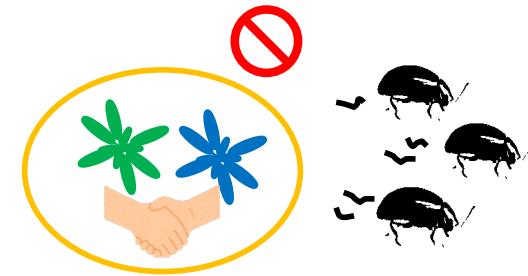


A researcher counts insects on the plants in the research garden on the UZH Irchel campus.

# Health of animals and plants depends on neighbors



Associational susceptibility



Associational resistance  
(positive biodiversity effect)



Having right or wrong neighbors?

Mixed stands are less or more damaged by herbivores (Jactel et al. 2021)

# *Arabidopsis* interact with other plants and natural enemies in the field

Lab (solitary)



Outdoor (group)



# Pest insect community on *Arabidopsis thaliana* in Irchel



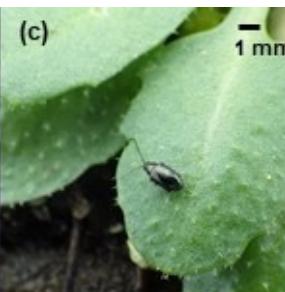
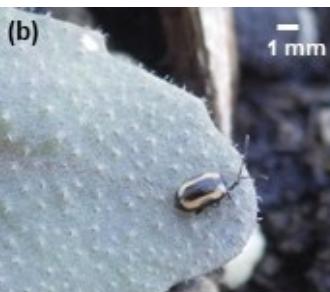
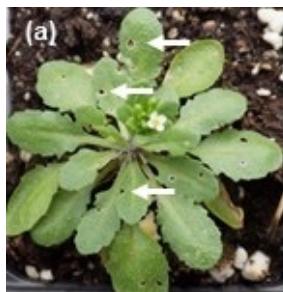
Zurich (natural range)      Shiga, Japan (naturalized range)

leaf holes made  
by flea beetles

a striped flea beetle  
*Phyllotreta striolata*

turnip flea beetle  
*Phyllotreta atra*

mustard aphids  
*Lipaphis erysimi*



turnip sawfly  
*Athalia rosae*

cabbage butterfly  
*Pieris rapae*

diamond back moth  
*Plutella xylostella*

western flower thrips  
*Frankliniella occidentalis*

# Aim: Genomic prediction of positive plant-plant interactions

*Arabidopsis thaliana*

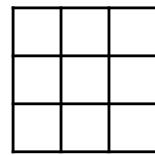


Insect herbivory

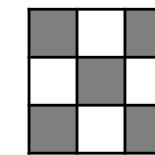


*Randomized mixture*

**(1) Picking up genotypes pairs that gain associational resistance**

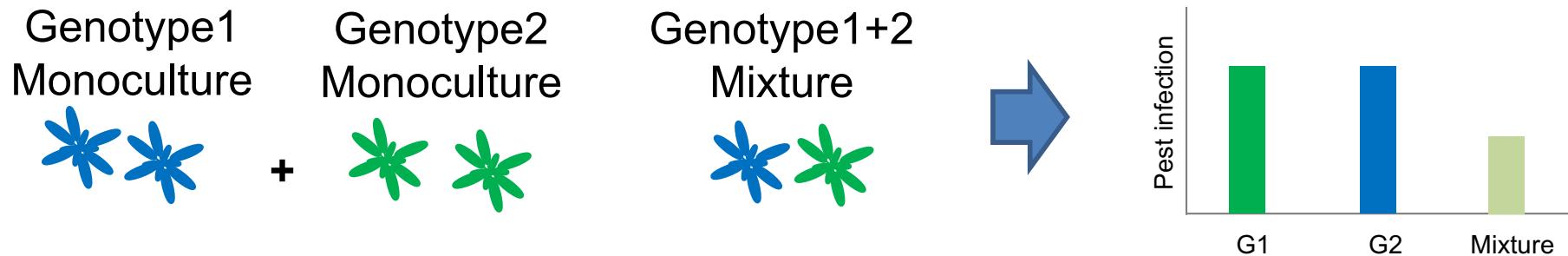


*Monoculture vs. Mixture*



**(2) Verifying key pairs with associational resistance**

# When does genetic diversity matter? Finding beneficial pairs for mixed planting (Wuest et al. 2021)



$${}_{200}C_2 = \frac{200*199}{2} = 19900 \text{ pairs from 200 genotypes...}$$

**Q. How can we find better pairs out of many possible pairs?**

# Ising model to study neighbor interactions

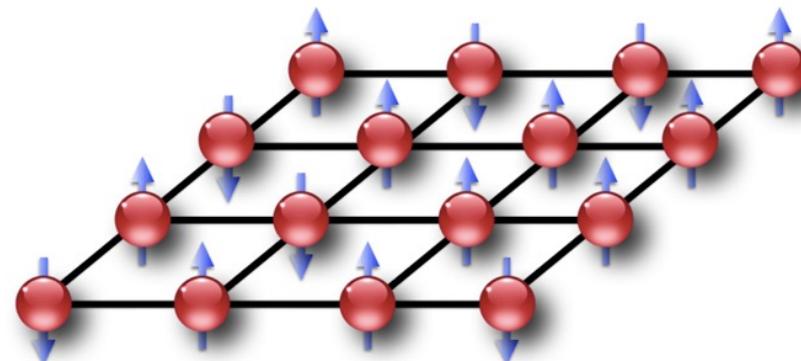


mathematical model of ferromagnetism  
in statistical mechanics



Ernst Ising (1900-1998)

Taroni, *Nature Physics* 2015



<https://nedo-quantum.aist.go.jp/isng-machine.html>

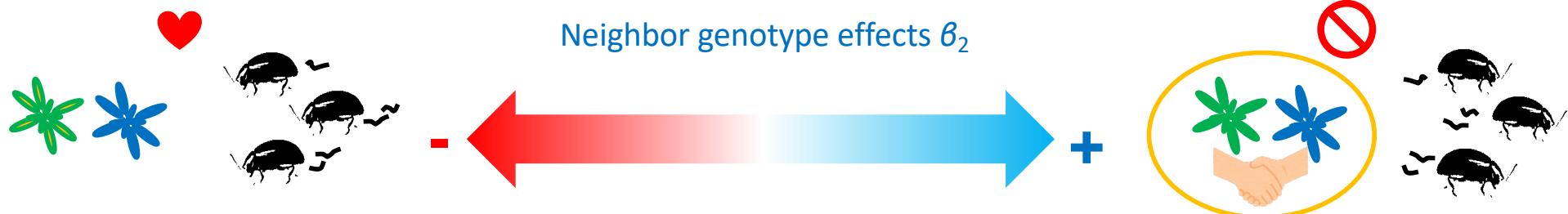
# New method Neighbor GWAS: Detecting neighbor interactions in randomized mixtures of genotypes

Normal GWAS: **Phenotype = Genotype + Environment**

Neighbor GWAS: **P = G + E + GxE (=GxNeighborG)** (Sato et al. 2021 *Heredity*)

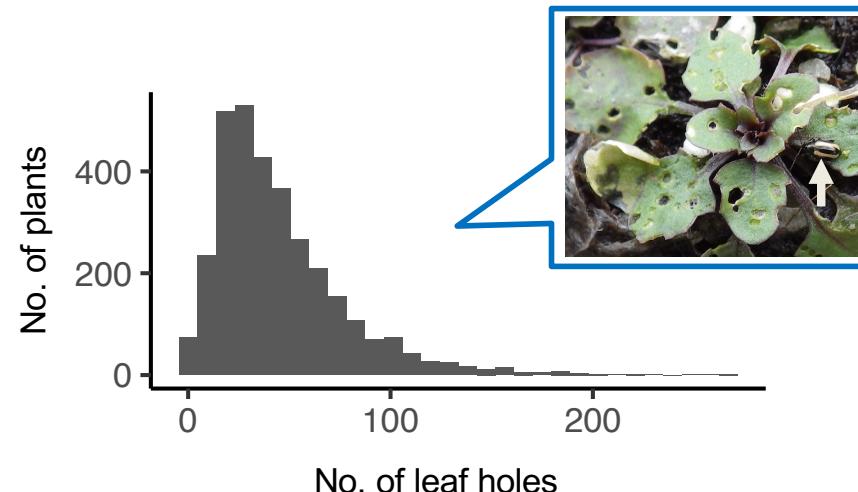
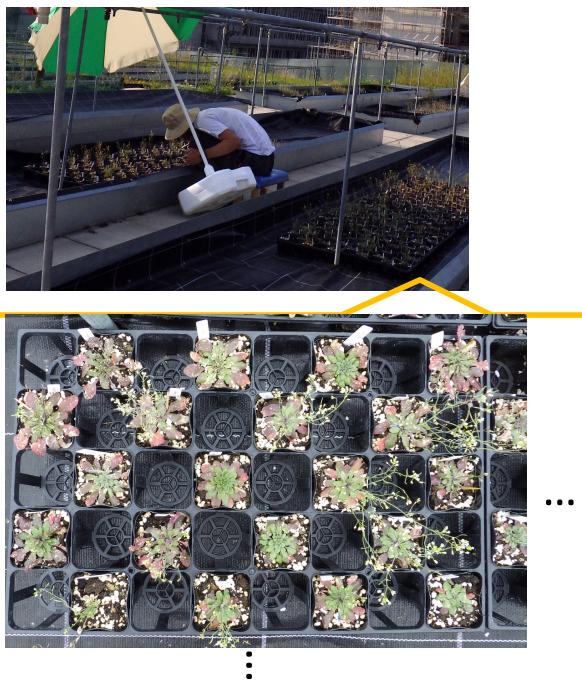
$$y_i = \beta_0 + \boxed{\beta_1 g_i} + e_i + \boxed{\beta_2 \frac{\sum_{j=1}^J g_i g_j}{J}}$$

focal G                            GxNeighborG



# Large dataset by field study in Irchel

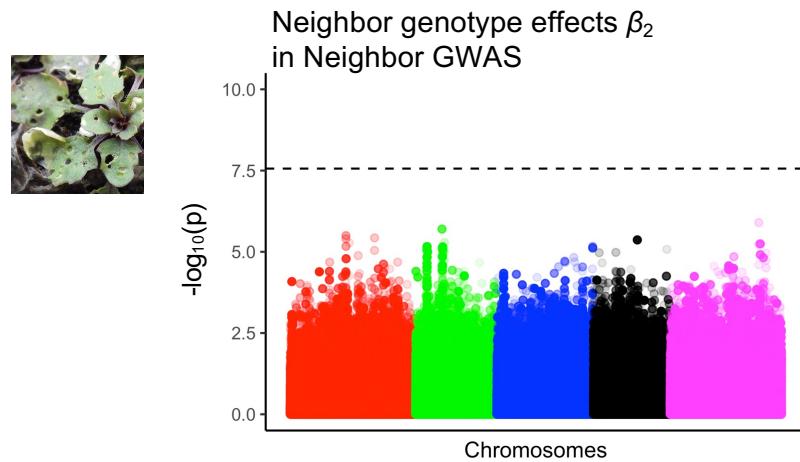
During July at Zurich, Switzerland



- 200 genotypes randomly assigned within a block
- × 8 blocks × 7 survey × 2 years × 2 sites = **44,800** data points

Observation of leaf damage, and 52,007 insects

# Genome-wide association studies (GWAS): no significant peak suggesting polygenic nature

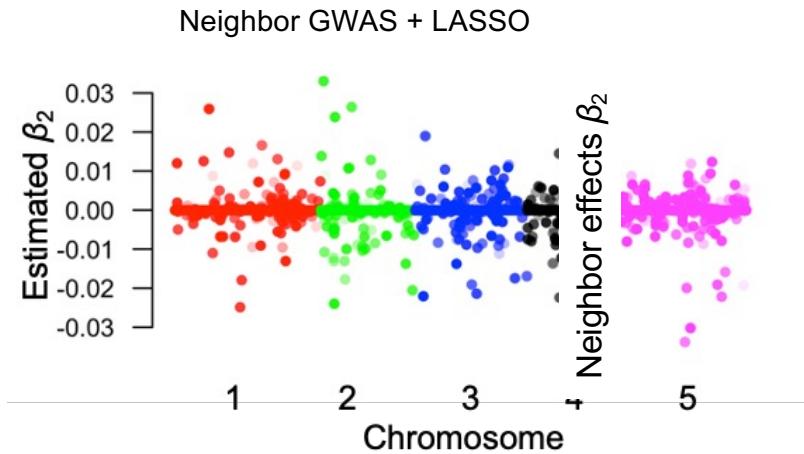


- Neighbor effects are unlikely attributed to a few loci and likely have a polygenic basis

# Machine learning (LASSO regression) both mechanisms and prediction

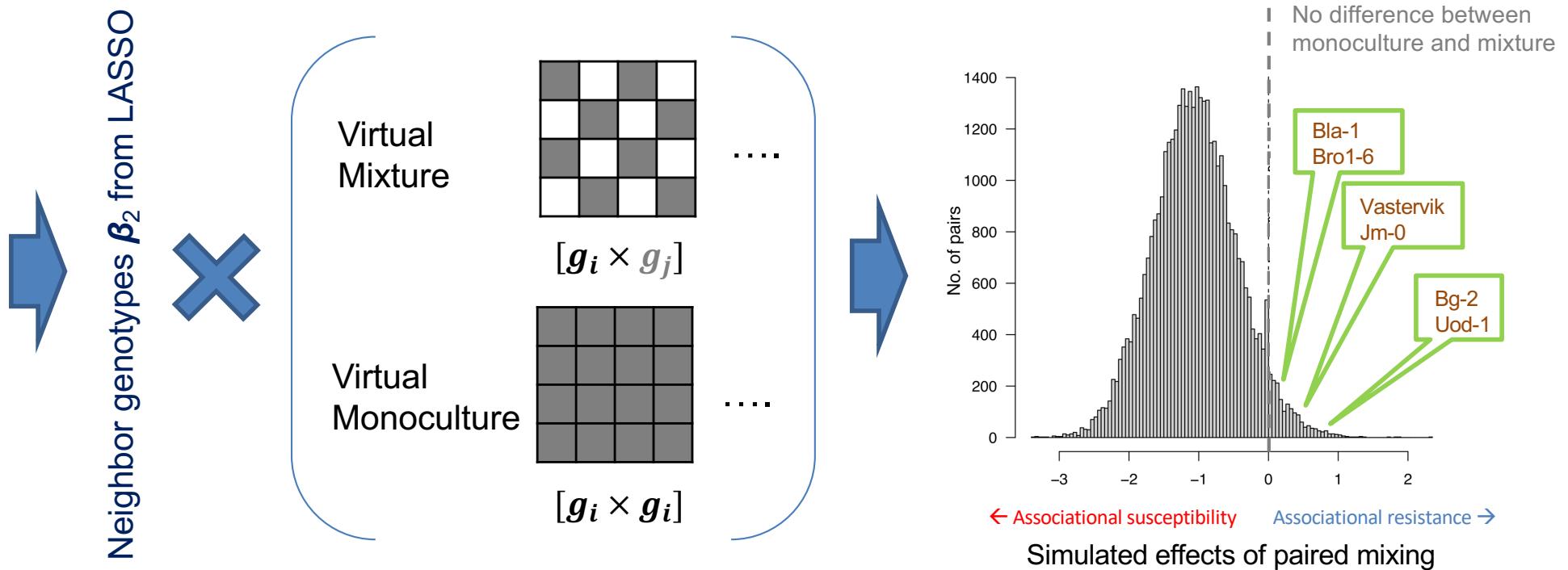
$$Y \sim a_1 * x_1 + a_2 * x_2 + \dots + a_n * x_n$$

Most of  $a_i$  to be zero

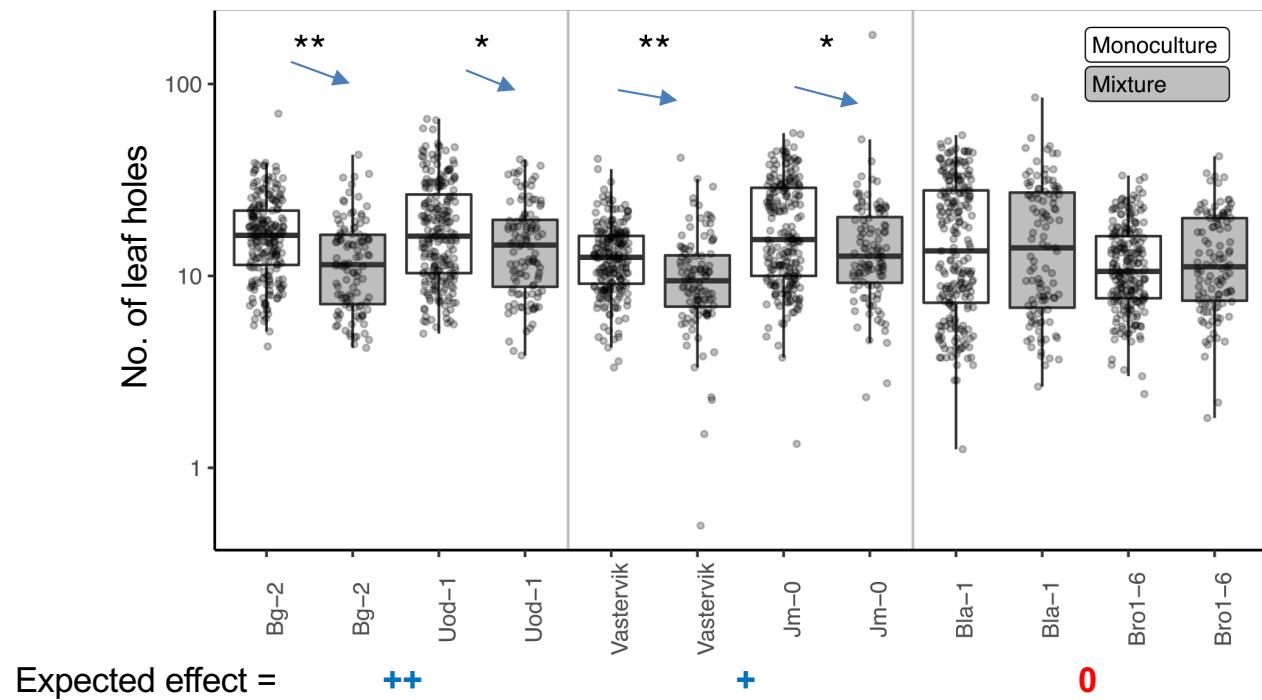
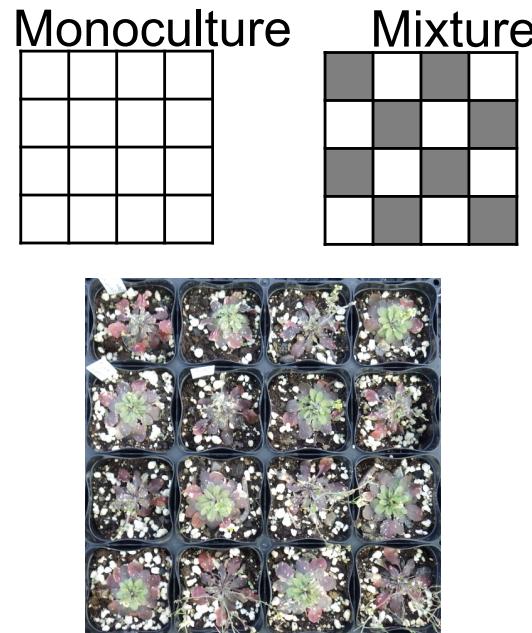


- LASSO regression narrowed down from 1.2 million to 756 SNPs

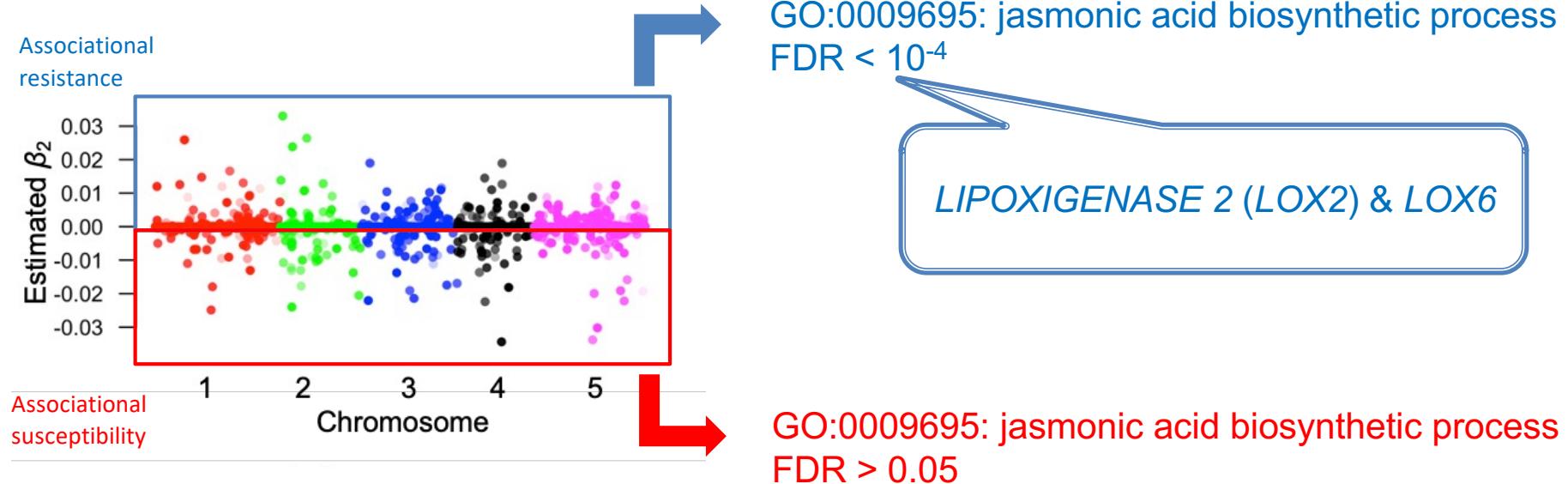
# Genomic prediction of herbivore damage under virtual mixture and monoculture



# Validating the prediction by LASSO regression: two genotype pairs indeed mitigated herbivory in mixed planting



# Potential mechanism? JA-related GOs are enriched in SNPs associated with associational resistance ( $\beta_2 > 0$ )



# Aim: Genomic prediction of positive plant-plant interactions

*Arabidopsis thaliana*



Insect herbivory



*Randomized mixture*

**(1) Picking up genotype pairs that gain associational resistance**

Predicting genotype pairs with positive biodiversity effects

**(2) Verifying key pairs with associational resistance**

# Insecticide: Food security and environment



Biodiversity



Food security



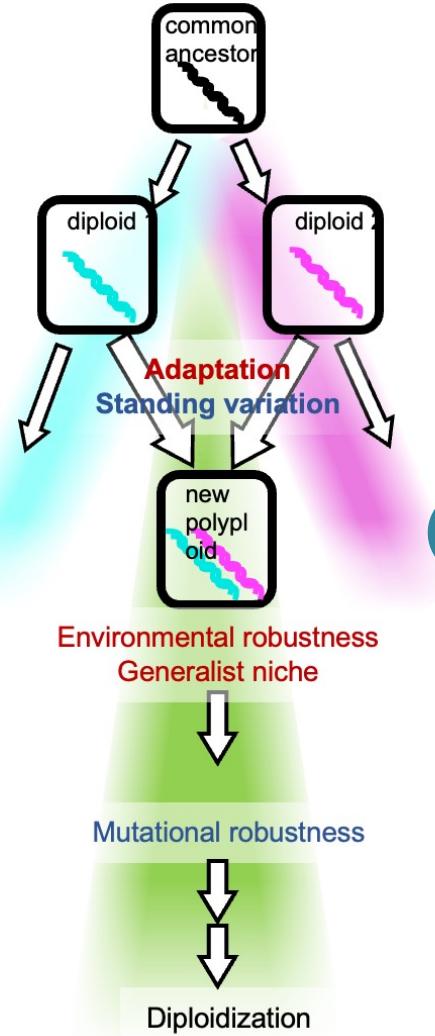
# Take-home message

1. Neighbor GWAS + LASSO found key genotype pairs out of randomized cultivation in field-grown plants  
→ Any plant genomes in randomized mixtures
  
2. This approach is useful even when a target trait is highly polygenic  
→ Any quantitative traits of ecological and agricultural interest

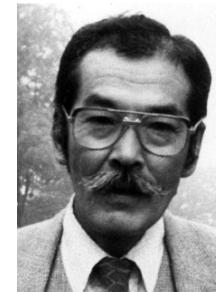
# Two published examples of genomic researches using machine learning

- Genome-wide association studies and genomic prediction of biodiversity effects
  - LASSO regression both for choosing interpretable variables and for prediction
- Hardware and software PlantServation for time-course image analysis of polyploid species in field
  - Deep learning, random forest

# Is genome duplication advantageous? Genomics and machine learning



"genome shock", "gigas" effect



Stebbins (1950, 1971)  
Polyplody may retard adaptation

Ohno (1970)  
Long-term  
evolutionary advantage  
'Evolution by  
gene duplication'

Why can polyploid species adapt to distinct or broader ecological niches?

# A textbook example of new species during past 150 years

Clonal propagation      Submergence tolerance

*C. rivularis* ( $2x = 16$ ; RR)



*C. amara* ( $2x = 16$ ; AA)



*C. insueta* ( $3x = 24$ ; RRA)



*Caradamine insueta*

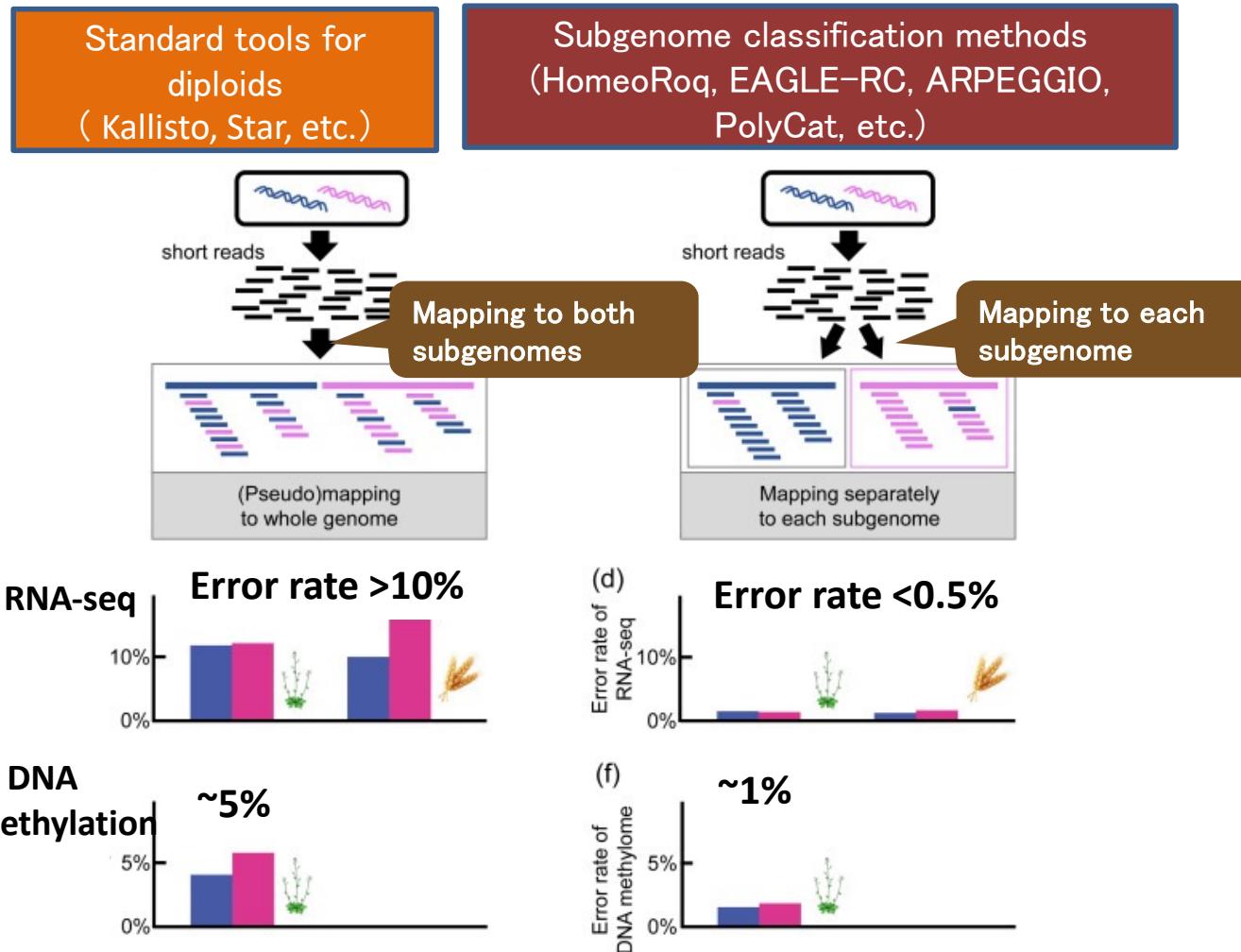


<https://www3.nhk.or.jp/nhkworld/en/tv/scienceview/>

Jianqiang Sun, Rie Shimizu-Inatsugi, Hugo Hofhuis, **Kentaro Shimizu**, Angela Hay, **Kentaro K. Shimizu**, Jun Sese, *Front Genet*, 2020

# Novel bioinformatic workflow for allopolyploid species

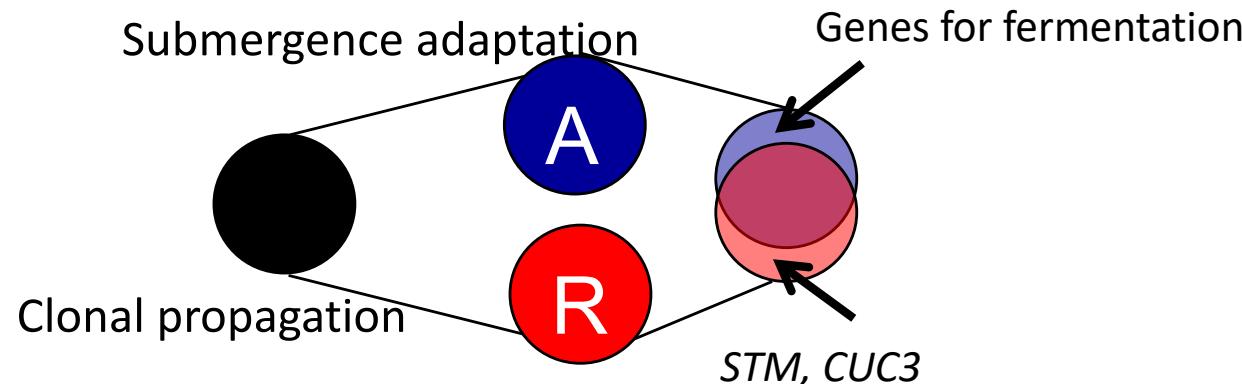
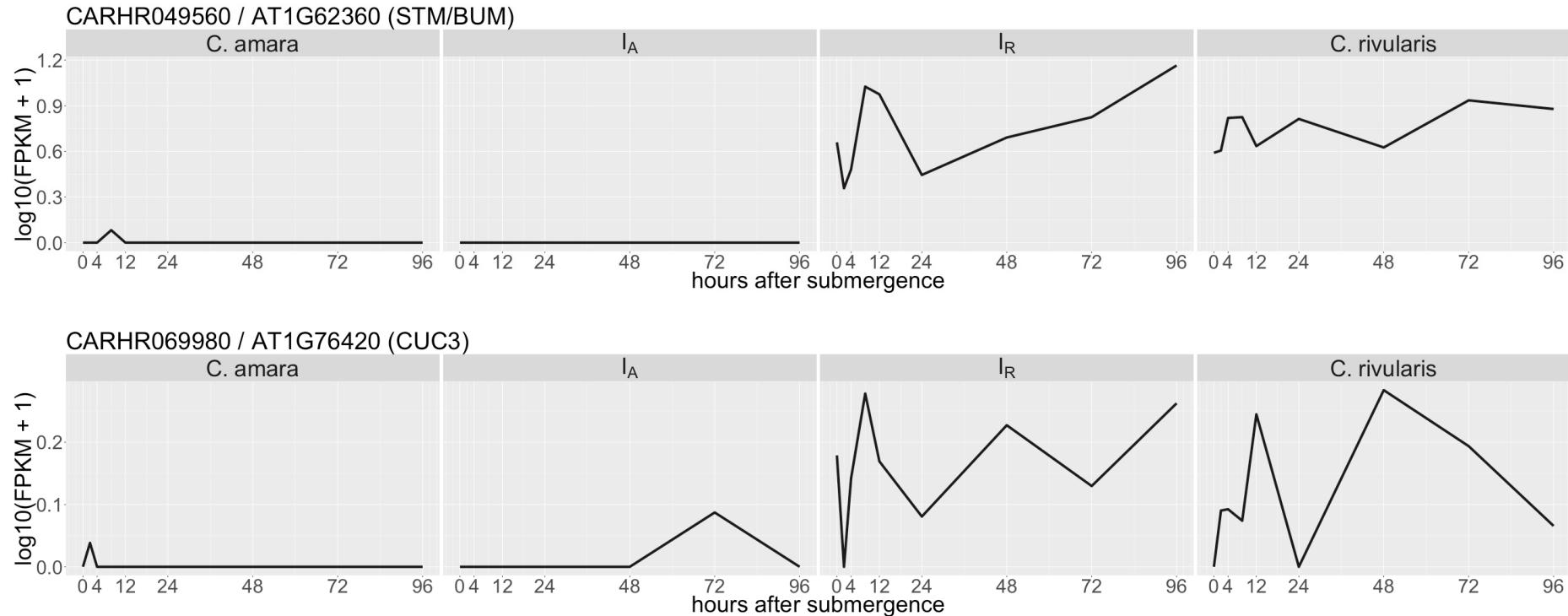
Shimizu, *Current Opinion in Plant Biology*, 69, 102292, 2022



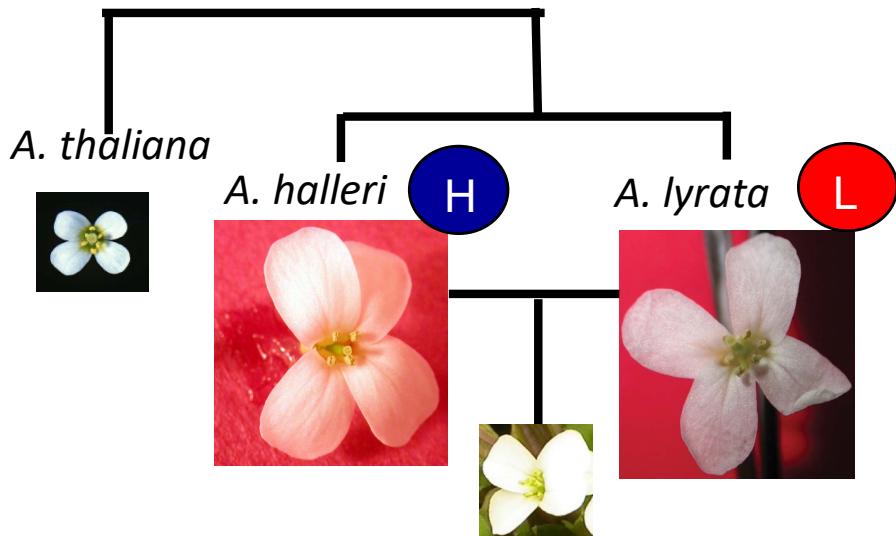
Akama *et al.*, NAR 2014, Kuo *et al.*, Brief Bioinf 2020, Milosavljevic *et al.* BMC Genomics 2022

Assembly, Resequencing, RNA-seq, Epigenetics

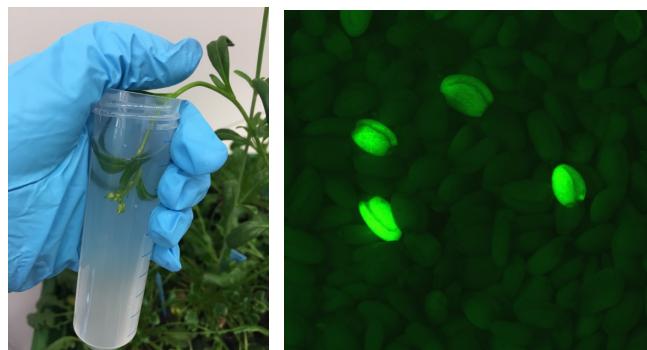
# Expression of key stem cell (meristem) gene was inherited



# Model polyploid species : the simplest for bioinformatic development: Synthetic and natural *Arabidopsis kamchatica*

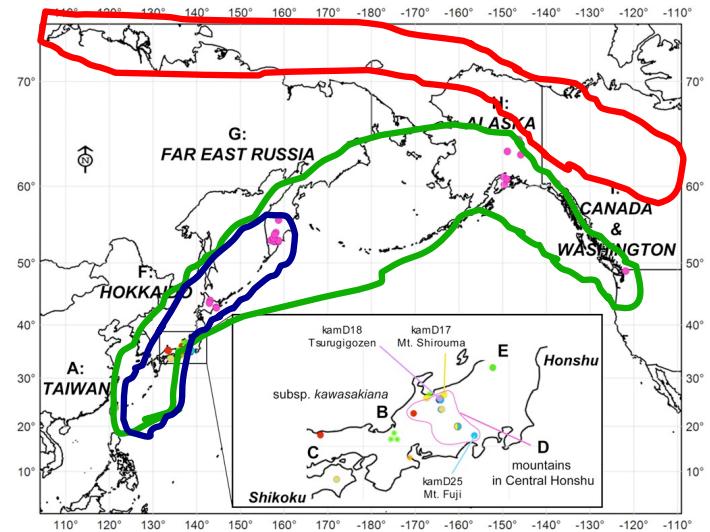


*A. kamchatica* (Fisch. ex DC.) K. Shimizu & Kudoh 2005



Transgenic technique  
Yew et al. *J Plant Res* 2018

Precipitation and temperature,  
Hoffmann, *Evolution* 2005

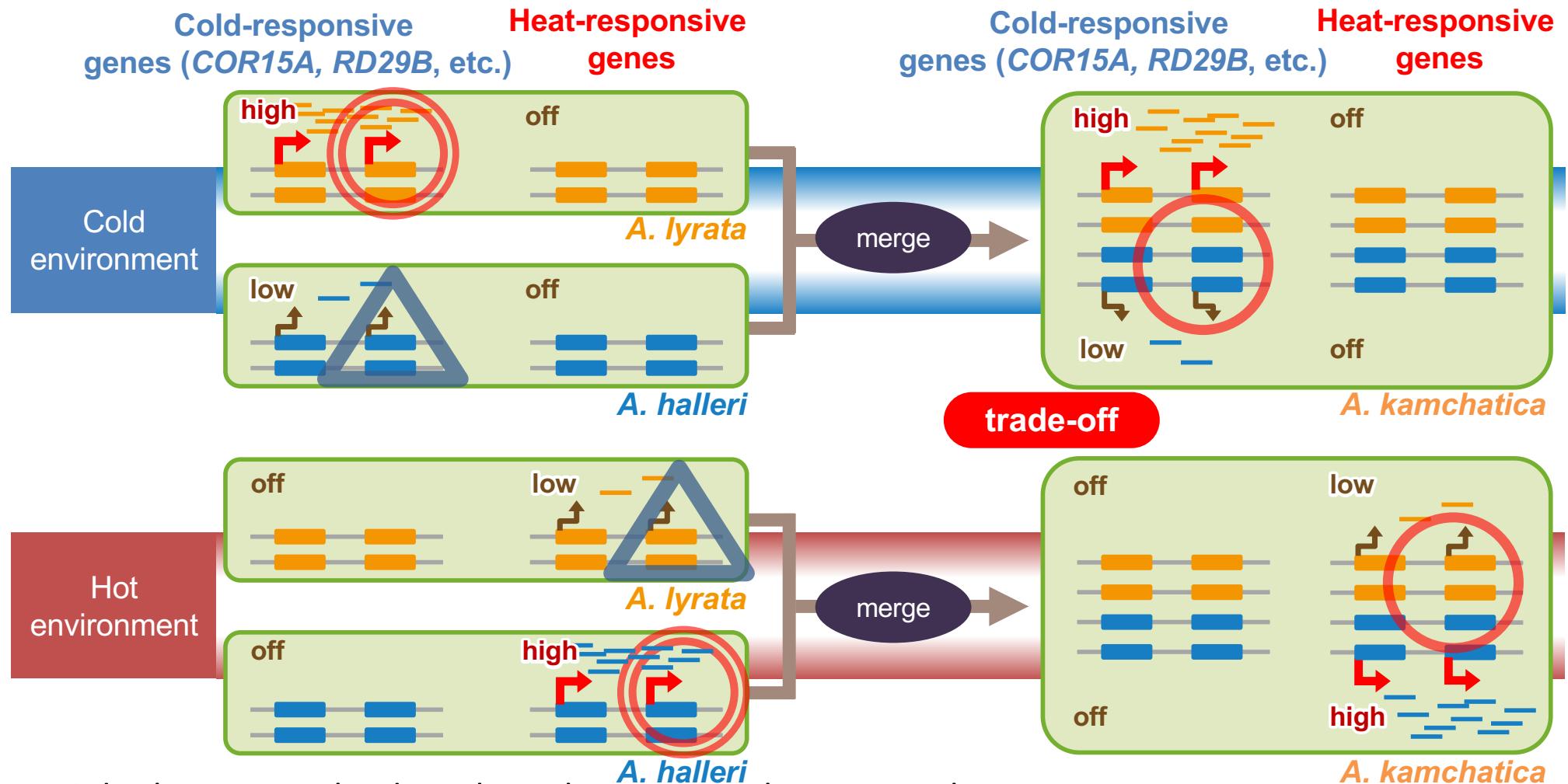


0-3000 m alt in Japan



Fujimoto et al. *PLoS Genet* 2008  
Shimizu-Inatsugi et al. *Mol Ecol* 2009  
Tsuchimatsu et al., *PLoS Genet* 2012

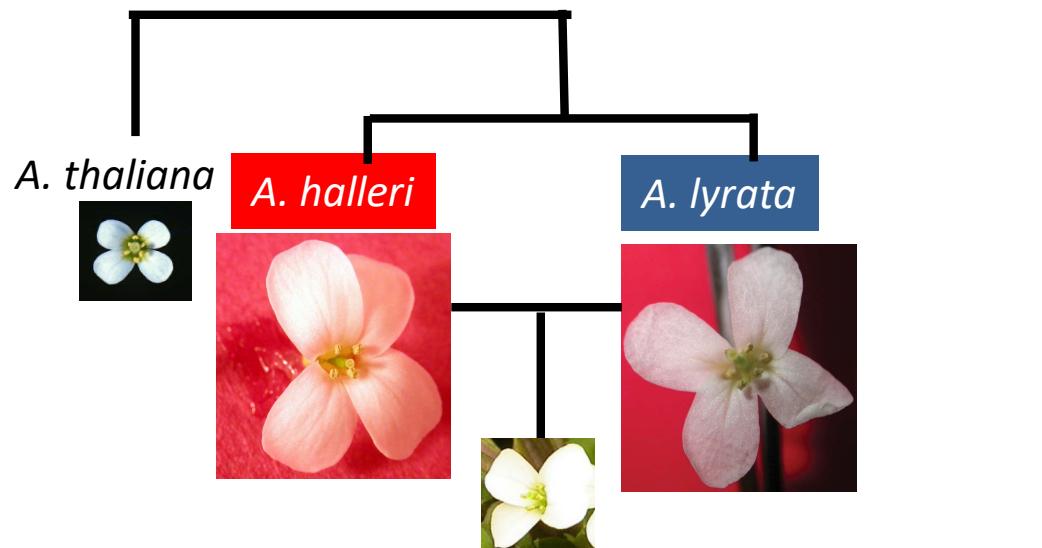
# Model of environmental robustness by inheriting and merging parental traits



Toward synthetic biology to confer robustness

# Phenomics using machine learning

## Often combined with genomics



*A. kamchatica* (Fisch. ex DC.) K. Shimizu & Kudoh

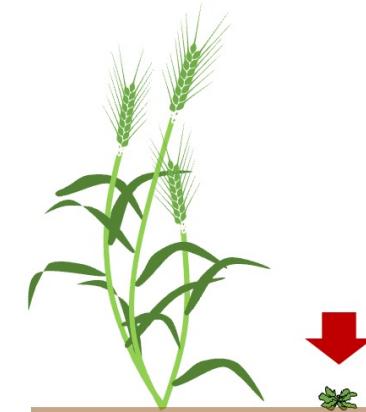
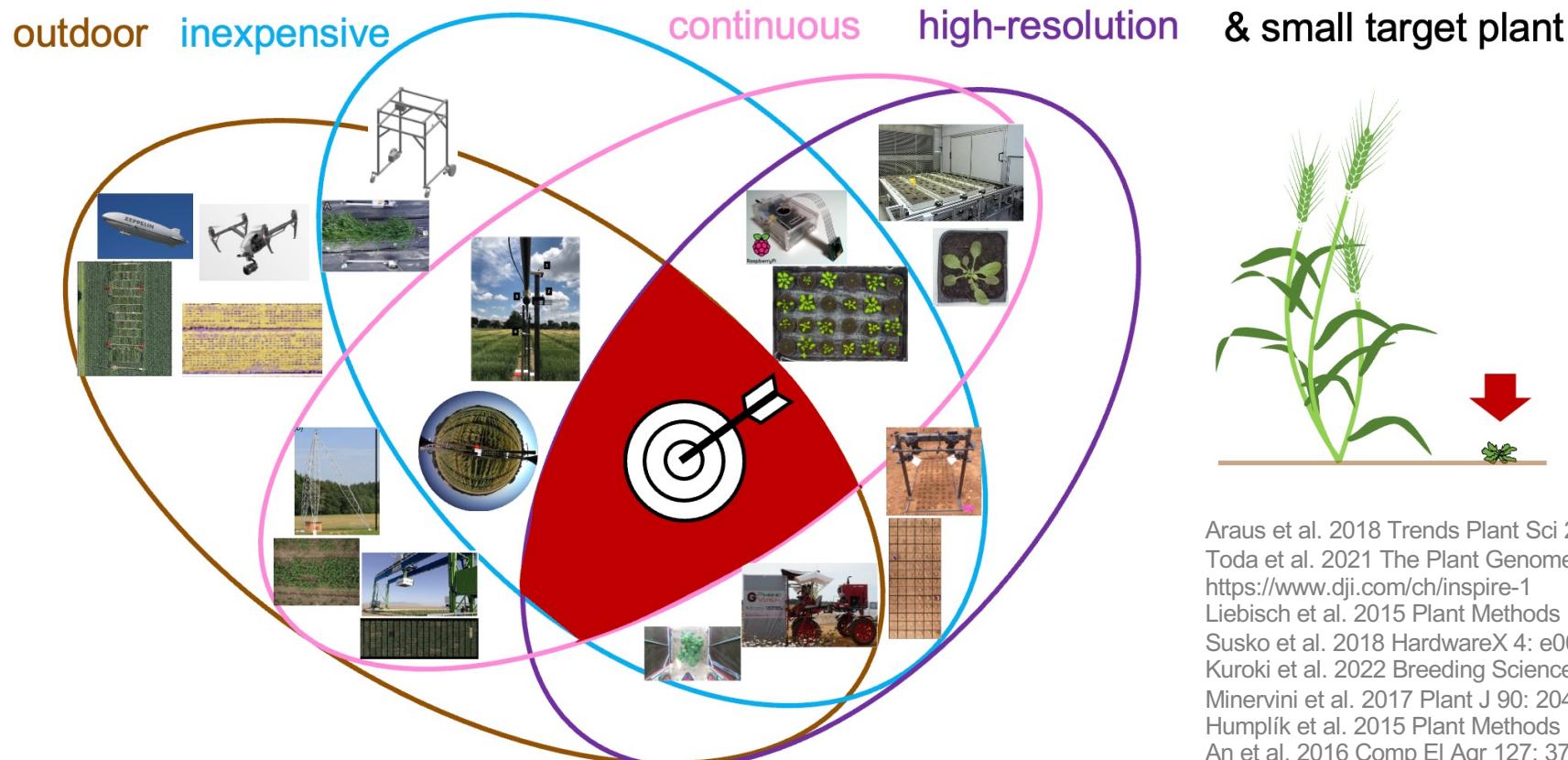


Stress marker of light, temperature, etc.

1. Did synthetic polyploids combine parental responses?
2. Can synthetic polyploids recapitulate natural speciation?

# Gap in monitoring hardware

Niche in remote sensing methods for accurate monitoring

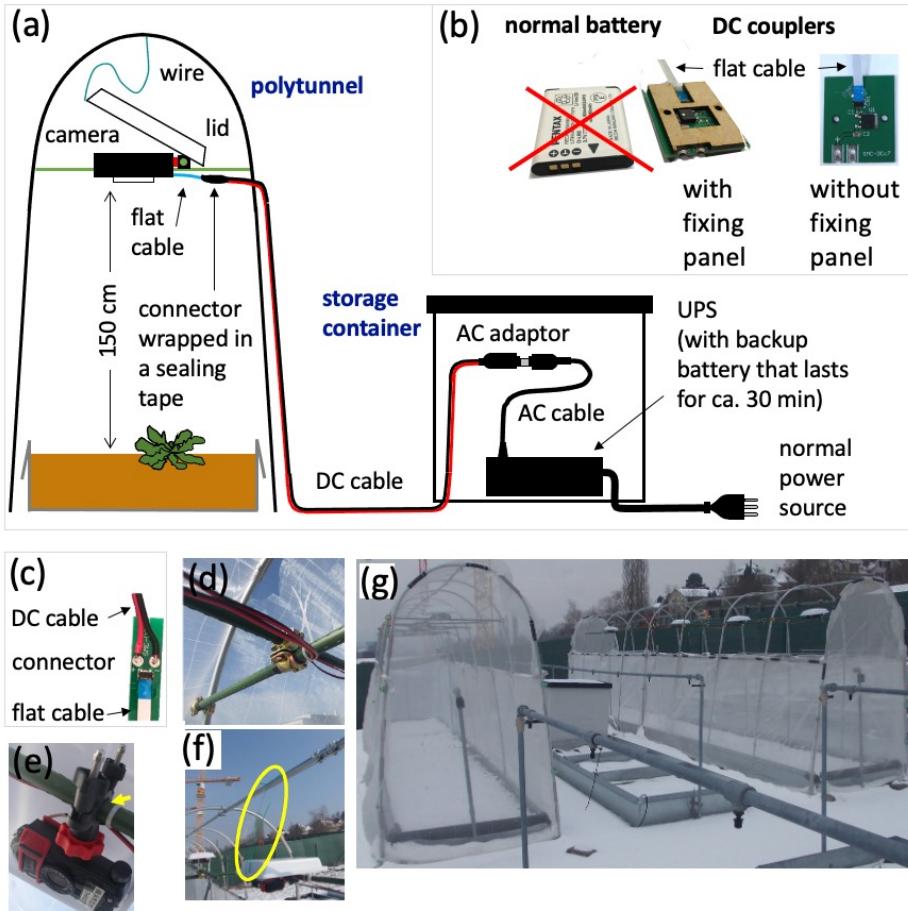


- Araus et al. 2018 Trends Plant Sci 23: 451-466  
Toda et al. 2021 The Plant Genome 14: e20157  
<https://www.dji.com/ch/inspire-1>  
Liebisch et al. 2015 Plant Methods 11: 9  
Susko et al. 2018 HardwareX 4: e00029  
Kuroki et al. 2022 Breeding Science 72: 66–74  
Minervini et al. 2017 Plant J 90: 204-216  
Humplík et al. 2015 Plant Methods 11: 29  
An et al. 2016 Comp El Agr 127: 376-394  
Jiang et al. 2018 Sci Rep 8: 1213  
Burnette et al. 2018 Proc Practice Experience Ad Res Comp 27: 1-7  
Kirchgessner et al. 2017 Func Plant Biol 44: 154-168

Figure by Reiko Akiyama

## New methodology of imaging analysis

# PlantServation: cost-efficient hardware and software



<http://www.nation.co.ke>

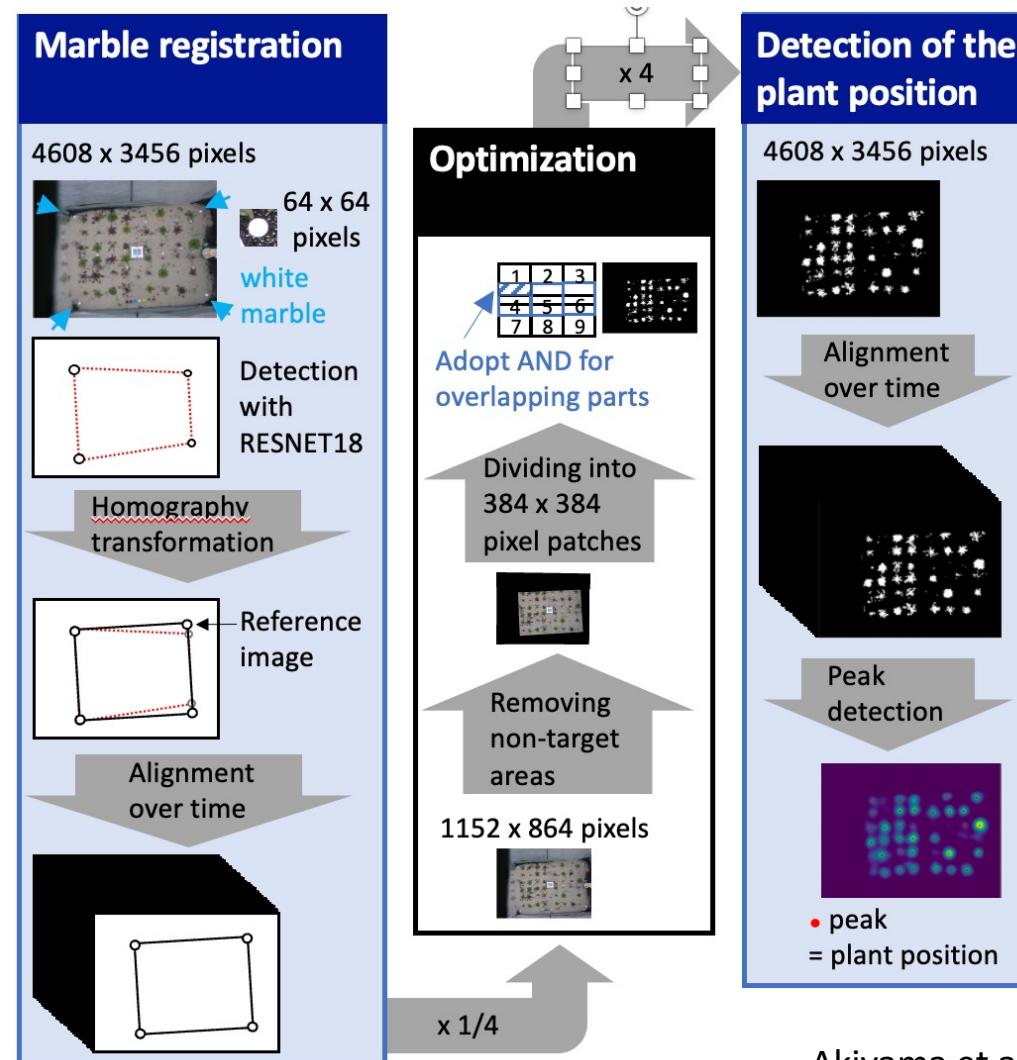


Akiyama et al., *Nature Commun* 14:5792, 2023

~2,600 USD for 384 individuals  
Suitable for researchers and small-scale cultivation

# Difficulties *in natura*: movement of camera

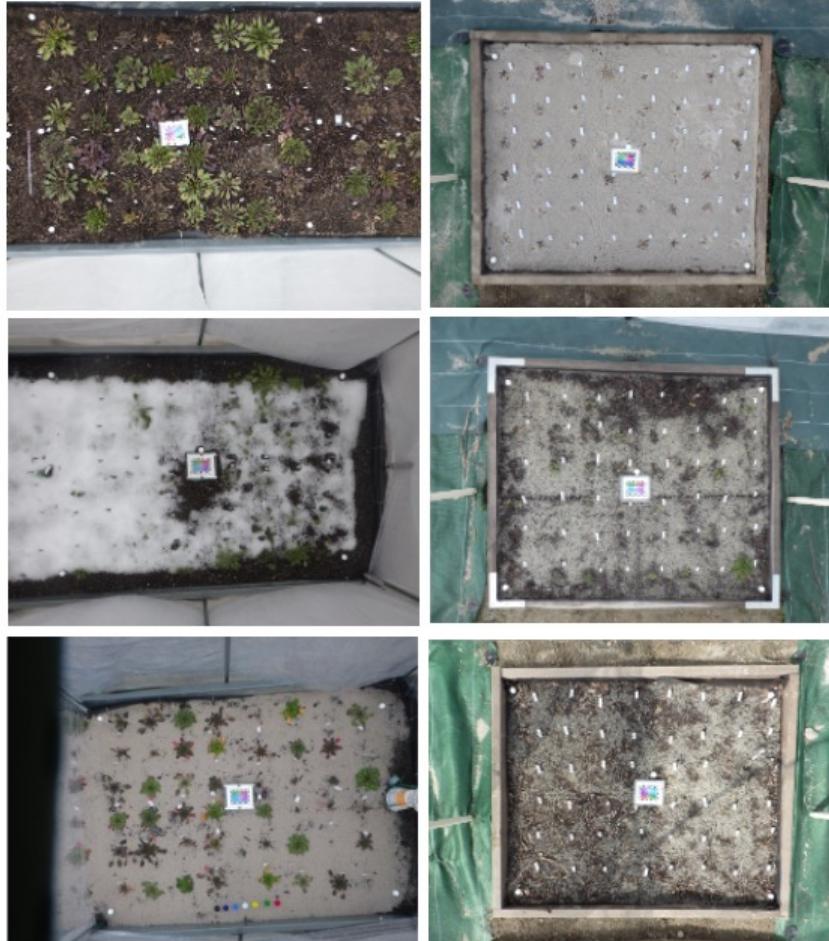
Wind, etc.



Akiyama et al., *Nature Commun* 14:5792, 2023

Identify plants using marbles

# Difficulties *in natura*: broad variations

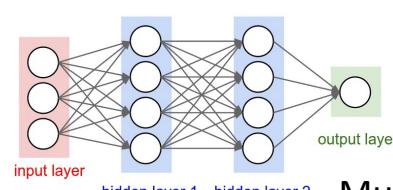
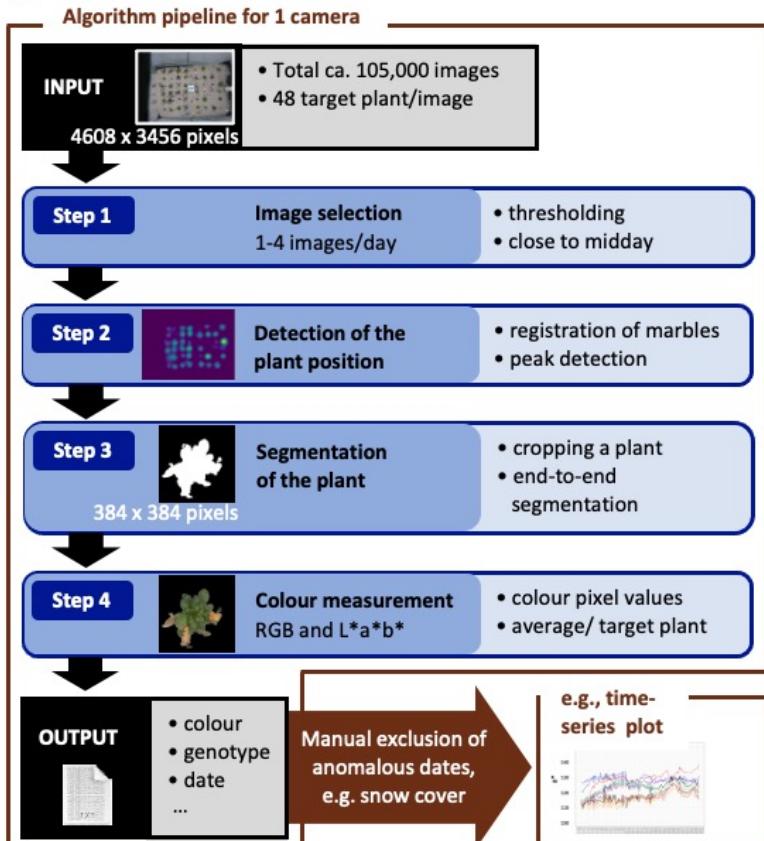


		broad leaf		thin petiole		thin leaf	
		red	green	red	green	red	green
sand + humus	small	sand				soil	
	dark	bright	dark	bright	dark	bright	dark

Akiyama et al., *Nature Commun* 14:5792, 2023

light condition and plants

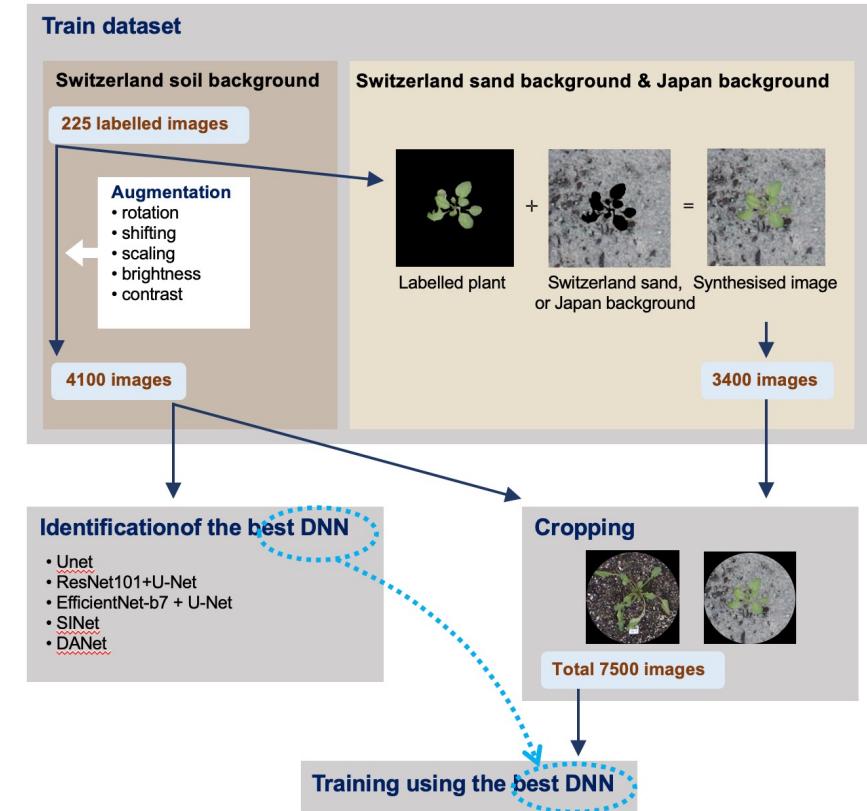
# Deep learning: manual labeling and augmentation



Musiol 2016



Manual labeling of 225 images (4 weeks)

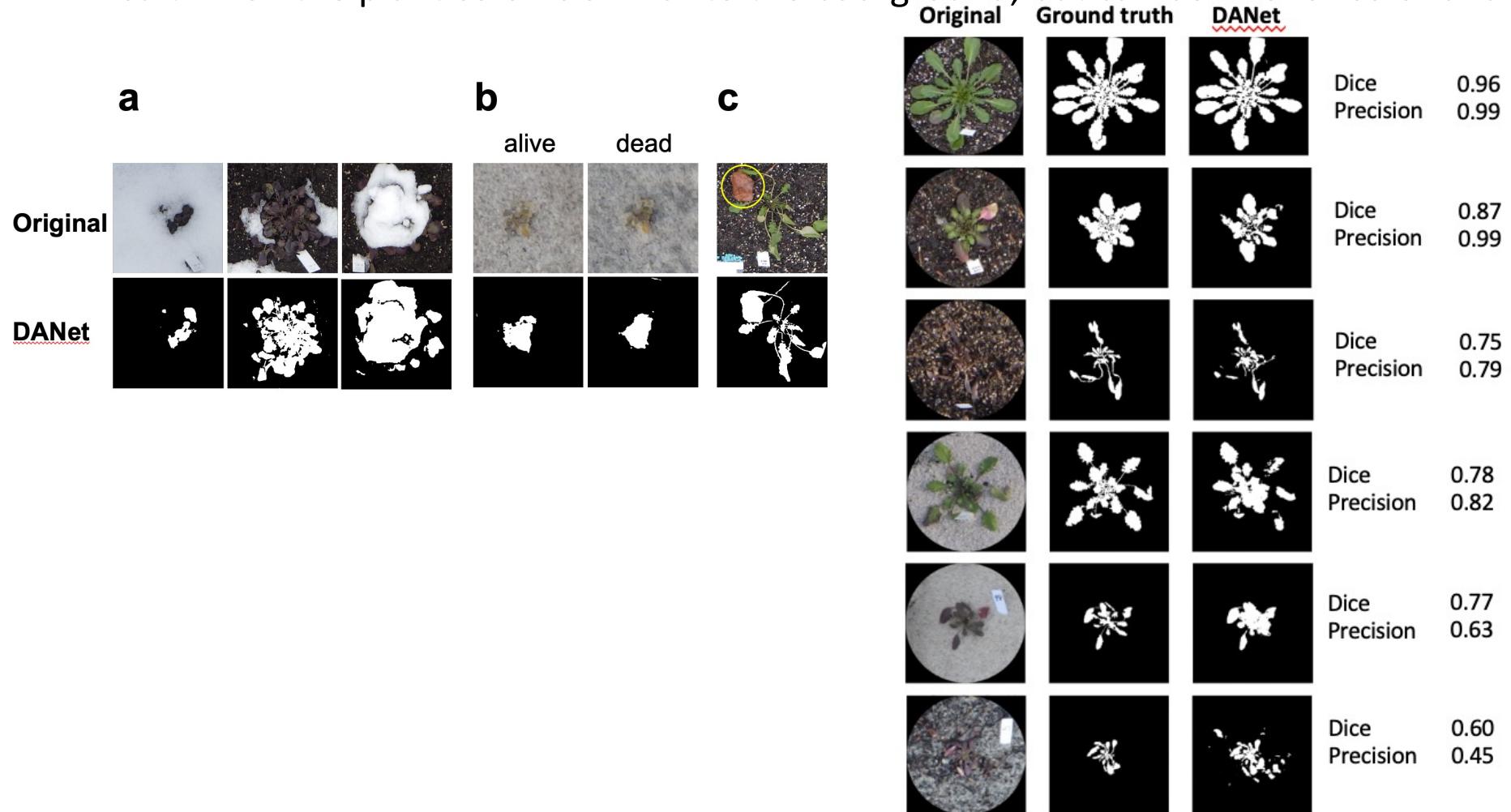


Akiyama et al., *Nature Commun* 14:5792, 2023

## Doable for biologists

# Checking by eye

Difficult when the plant color is similar to the background, but can be fine for color analysis



Akiyama et al., *Nature Commun* 14:5792, 2023

Remove snow season and outliers

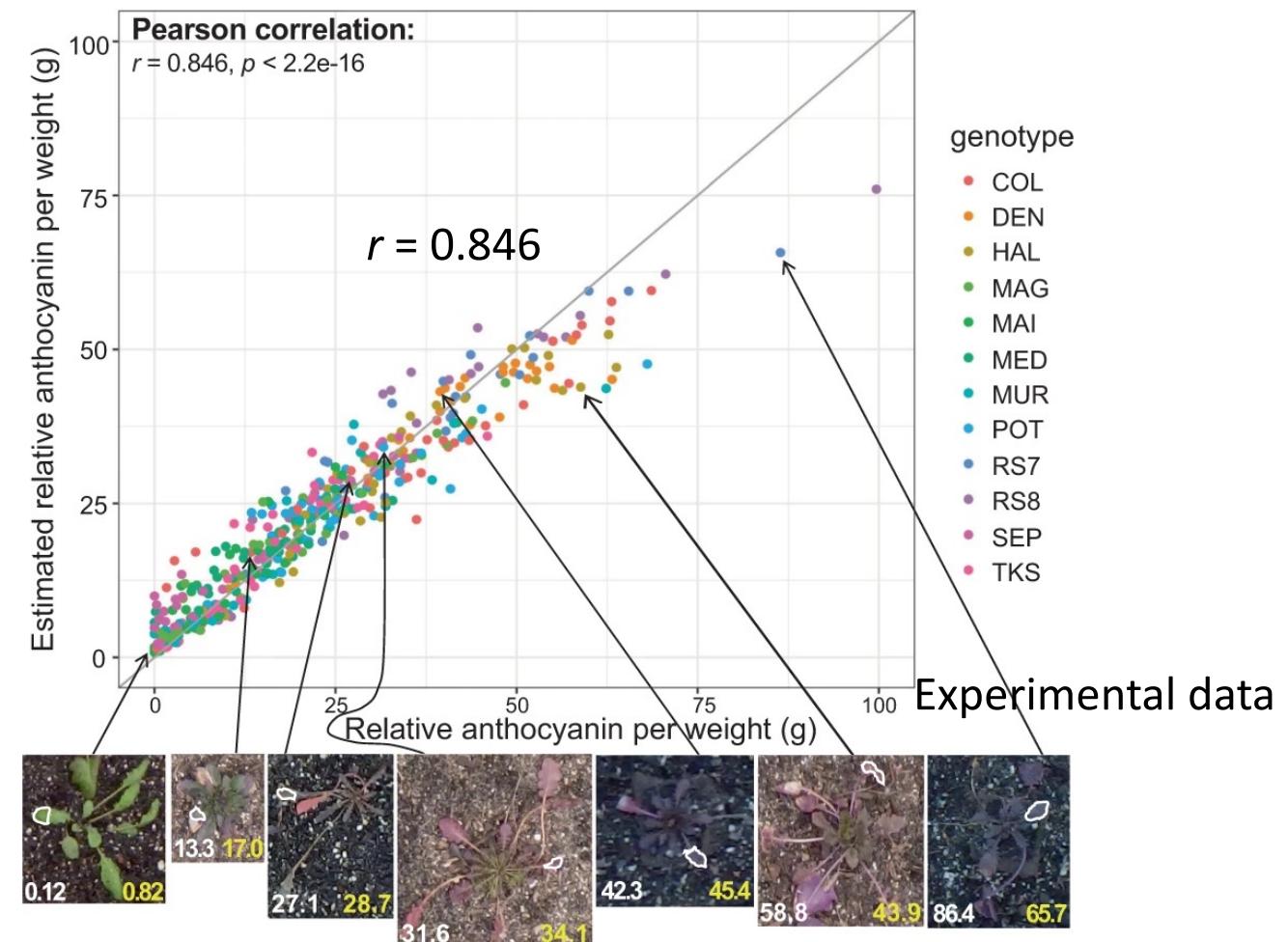
# Estimation and validation on anthocyanin concentration

Experimental (destructive)  
anthocyanin data of  
451 samples

Fitting

Random forest  
(a machine learning method)  
used for lab samples by Askey  
et al. 2019

App Plant Sci 7: 11 e11301



Akiyama et al., *Nature Commun* 14:5792, 2023

## Automatic pipeline established

# Machine learning: shape and color

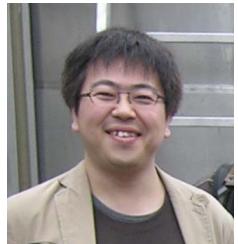
## Deep learning for segmentation



Manual labeling  
of 225 images  
(4 weeks)



Natsumaro  
Kutsuna,  
LPixel



Jun Sese  
Humanome

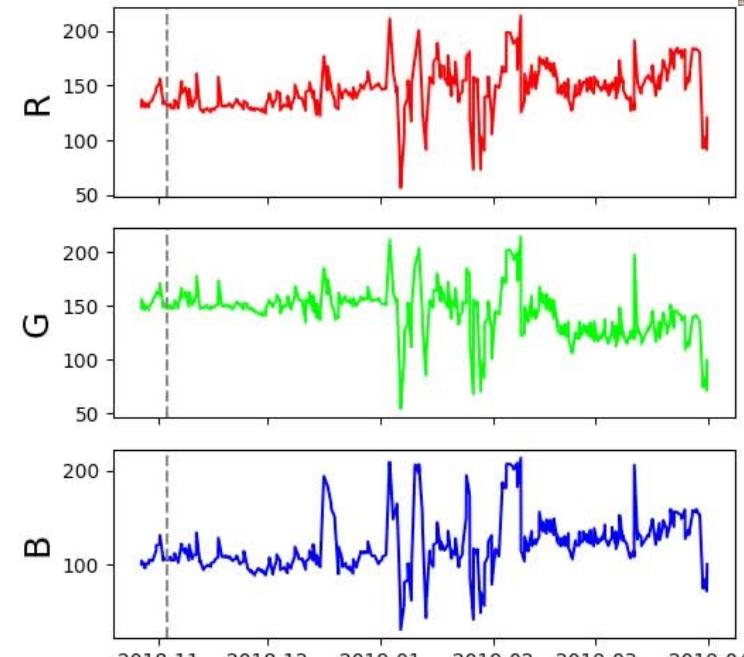


## Random forest for anthocyanin content

Experimental  
Anthocyanin data of  
451 samples



58.8      43.9

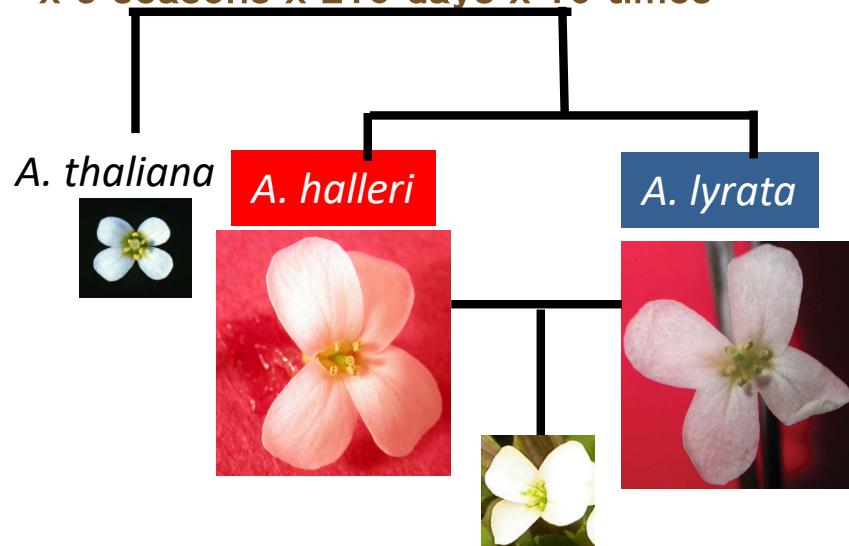


Akiyama et al., *Nature Commun* 14:5792, 2023

PlantServation pipeline established

# Analysis of 3,870,000 images using PlantServation

12 genotypes x 16 replicated x 2 sites  
x 3 seasons x 210 days x 16 times

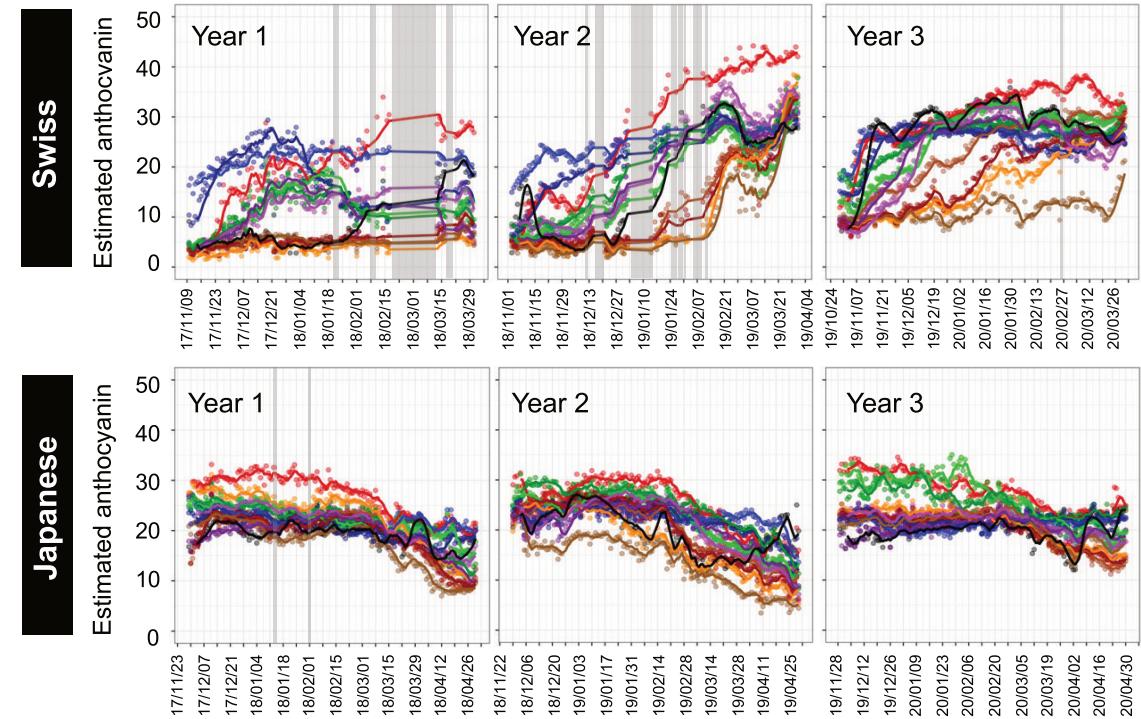


*A. kamchatica*  
2 synthetic polyploids  
2 natural (North)  
4 natural (South)

Zurich, Switzerland   Kyoto University, Japan



Hiroshi Kudoh

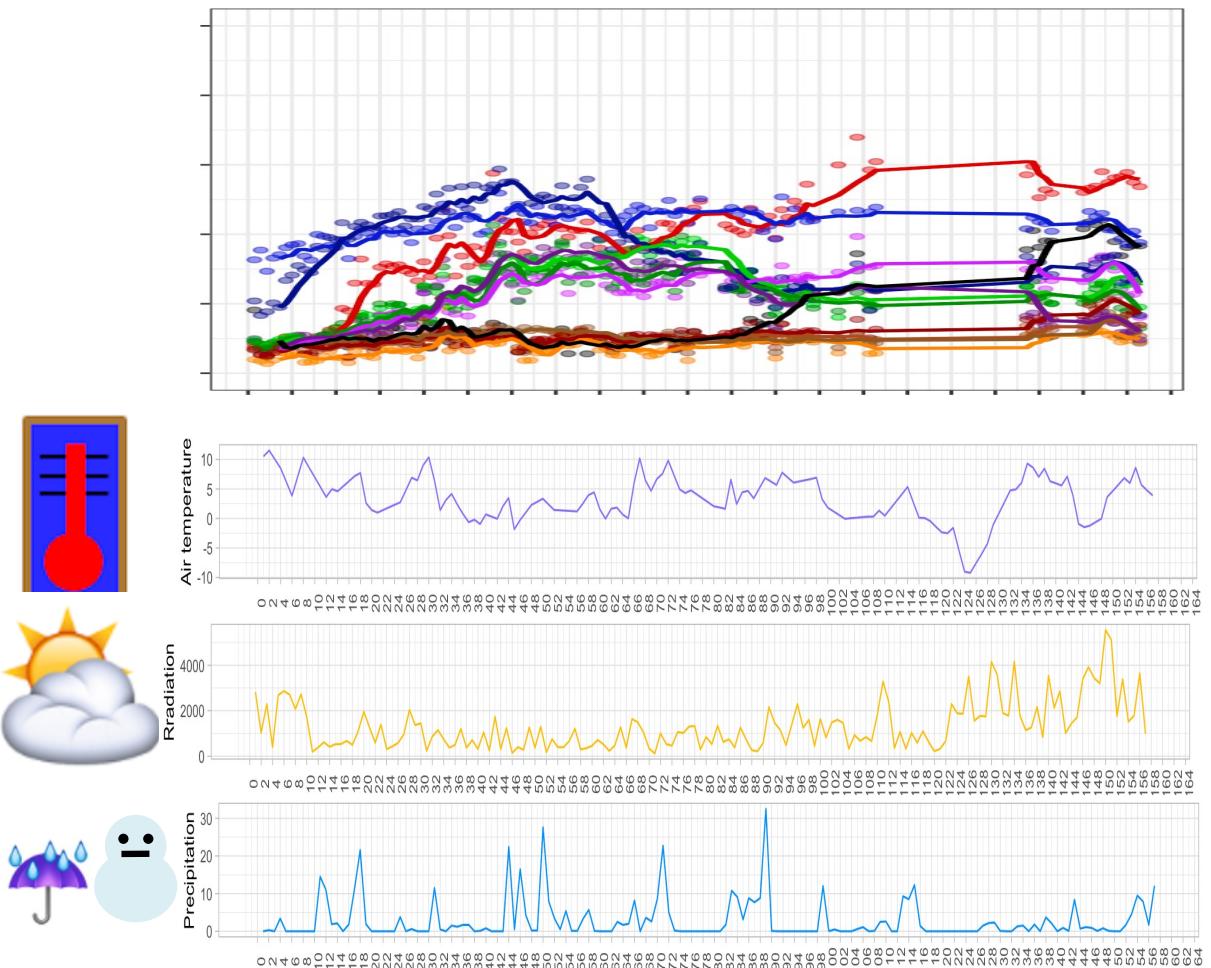


Akiyama et al., *Nature Commun* 14:5792, 2023

## High variation among genotypes and seasons

# Significant effect of temperature, sunlight, rainfall

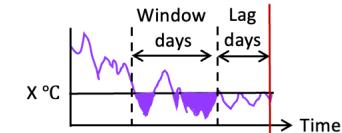
Zurich  
2017-2018



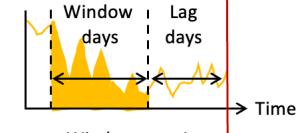
Genotype

- HAL
- MED
- SEP
- RS7
- RS8
- MAI
- TKS
- MAG
- MUR
- DEN
- POT
- COL

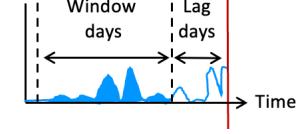
~5-10%



~5-20%



~5-30%

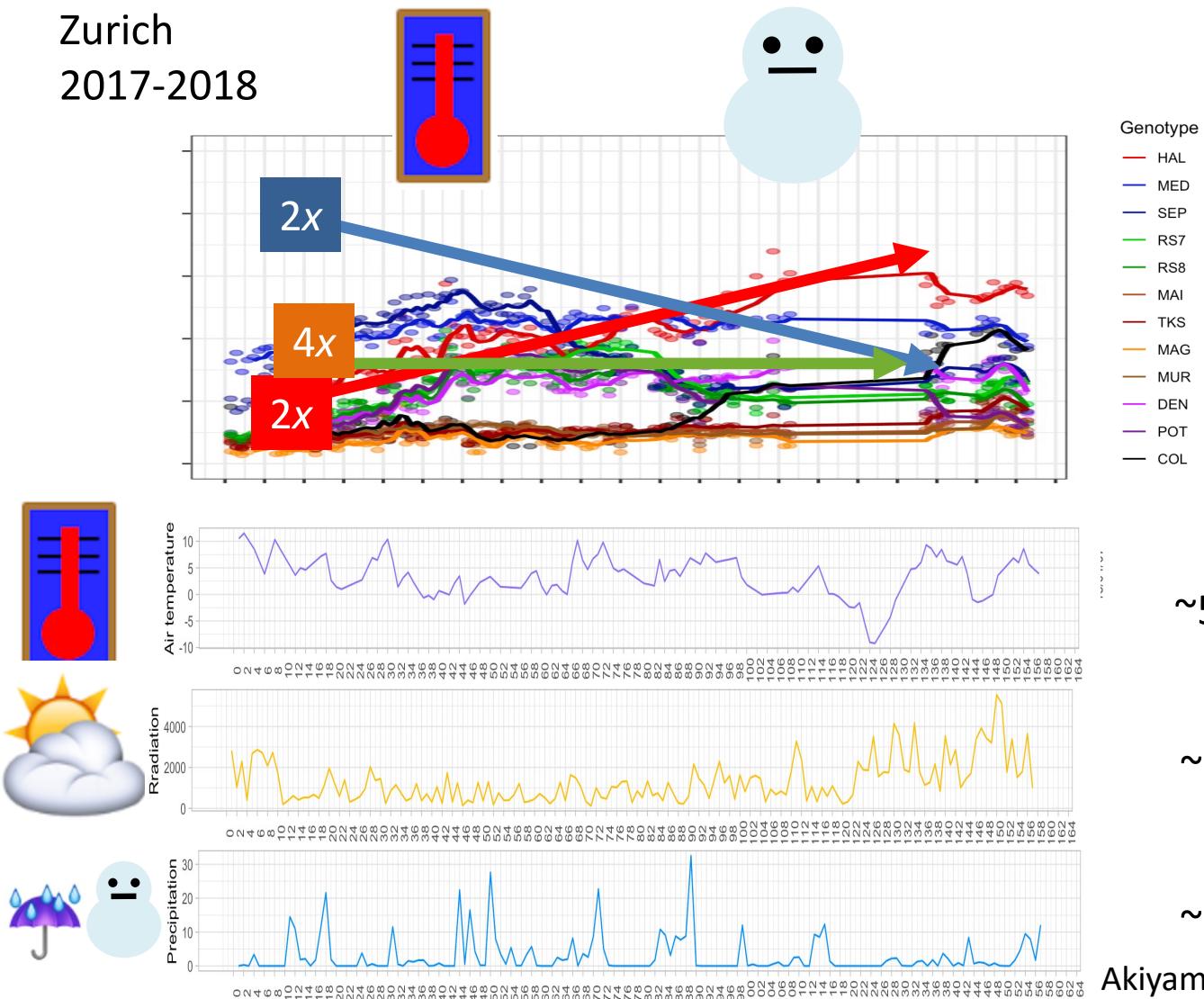


Akiyama et al., *Nature Commun* 14:5792, 2023

**Majority of variation unexplained:**  
**We know little about the environmental responses *in natura***

# Genotypic differences

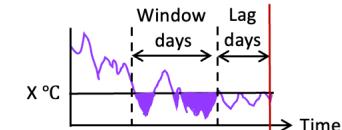
Zurich  
2017-2018



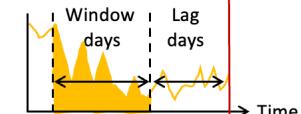
Genotype

- HAL
- MED
- SEP
- RS7
- RS8
- MAI
- TKS
- MAG
- MUR
- DEN
- POT
- COL

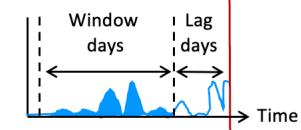
~5-10%



~5-20%



~5-30%

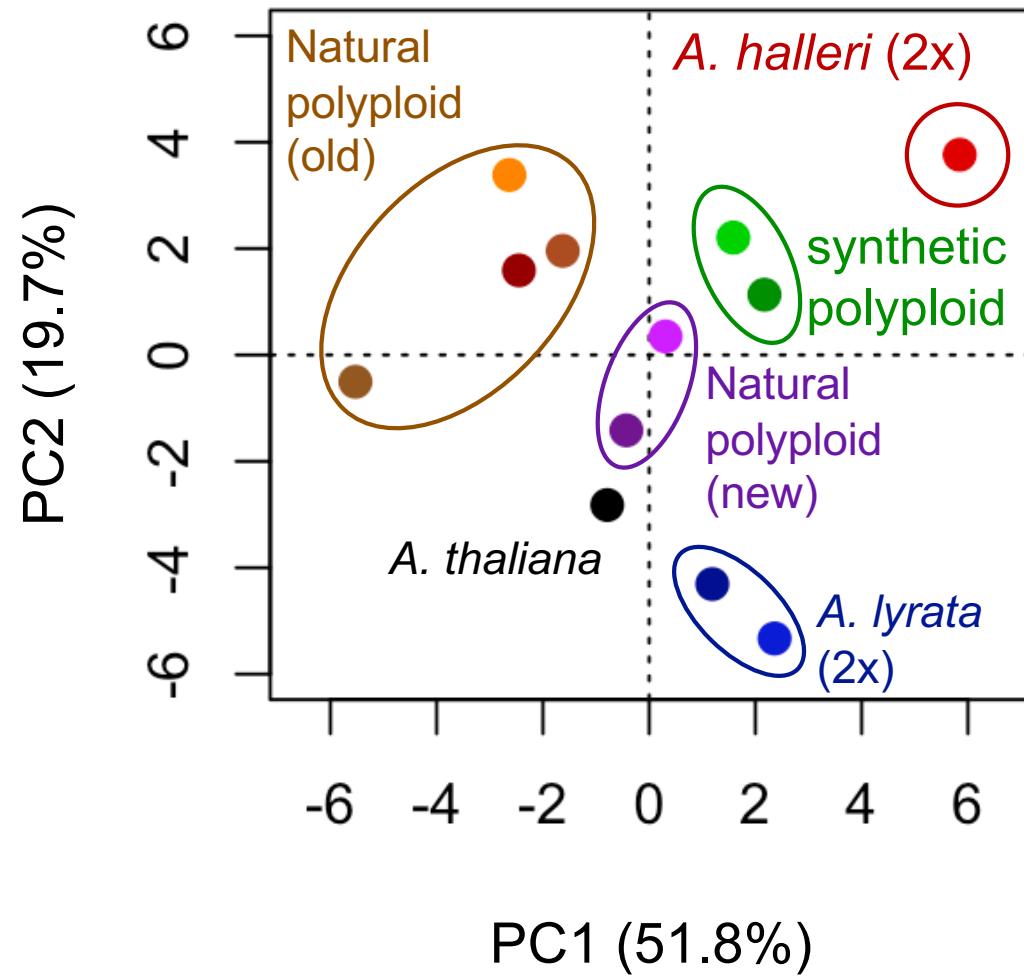
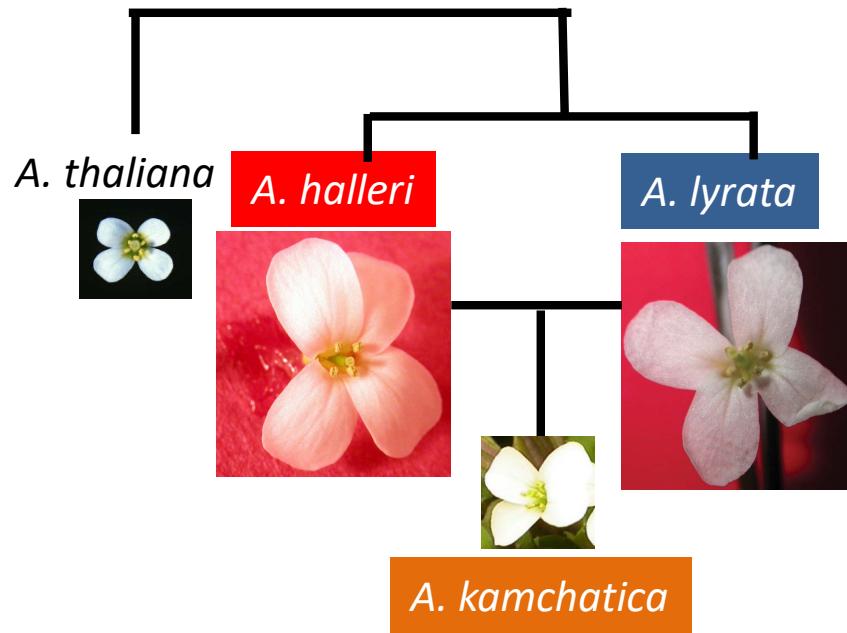


Akiyama et al., *Nature Commun* 14:5792, 2023

Polyploids were less stressed

# Synthetic polyploids combined parental responses and recapitulated natural ones

Principal Component Analysis

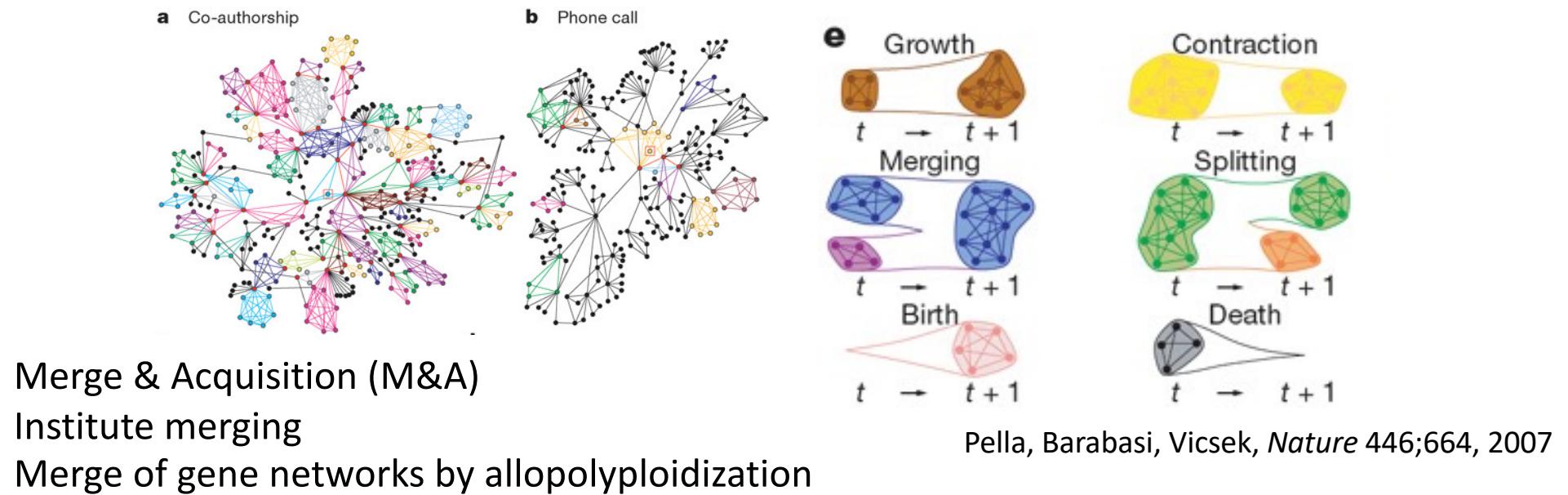


Akiyama et al., *Nature Commun* 14:5792, 2023

Did synthetic polyploids combine parental responses? YES

Can synthetic polyploids recapitulate natural speciation? YES

# Discussion: generalist by a network merging



Diploids	Allopolyploids
Specialist	Generalist
Expert (geek)	Management (jack-of-all-trades and mater of none)
Niche industry	conglomerate

Shimizu, *Current Opinion in Plant Biology*, 69, 102292, 2022

**Coexisting but agroecosystems may favor polyploids**