

# BIO334 Practical Bioinformatics

## Introduction to genome-wide polymorphisms

Take BIO373 for biological implications,  
or refer to the textbooks

1. Futuyma “*Evolution*”, Sinauer, 2005, 2017, Chap 12
2. Hartl and Clark “*Principles of Population Genetics*”, Sinauer  
4th ed., 2007
3. Freeman and Herron “*Evolutionary Analysis*” 5th Ed., Pearson, 2013,  
Chap. 7, 8

Department of Evolutionary Biology and Environmental Studies  
University of Zurich

Yasuhiro (Yasu) Sato, a senior postdoc of Shimizu Group

[yasuhiro.sato@uzh.ch](mailto:yasuhiro.sato@uzh.ch)

# Overview molecular population genetics

Introduction:

Genome-wide polymorphism and next-generation sequencers

1. nucleotide diversity

selective sweep / hitchhiking, and balancing selection

2. Tajima's  $D$

3. positive selection on replacement mutations

Human 1000 genomes were published on 28 Oct 2010:  
How are they meaningful?



# Example: evolution of language

Was there positive selection to have language?

How many genes were necessary to obtain language?

How and when did the language evolve?

Did Neanderthal have language,  
which diverged 3-400,000 years ago?

Limited evidence: archaeology, fossil of skeletons

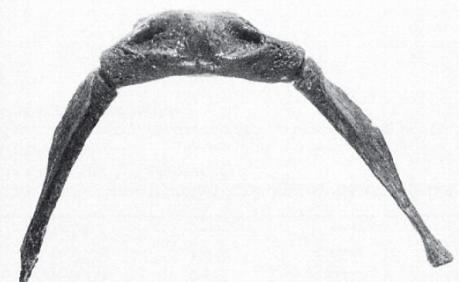
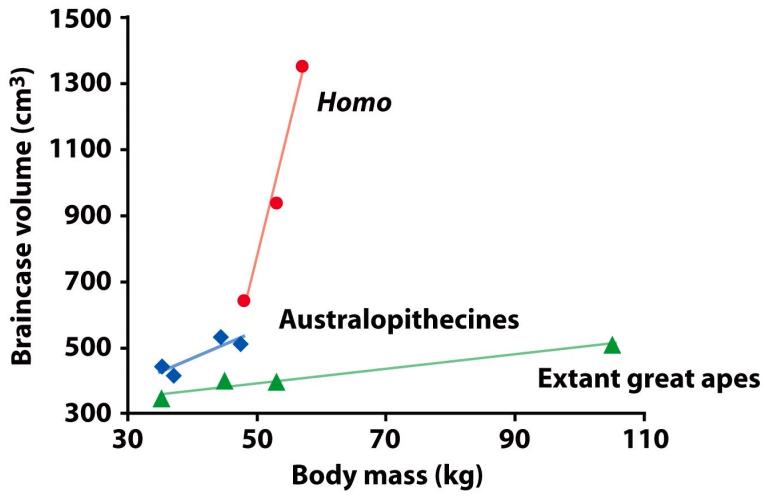


Figure 20-31 Evolutionary Analysis, 4/e  
© 2007 Pearson Prentice Hall, Inc.

Figure 20-32 Evolutionary Analysis, 4/e  
© 2007 Pearson Prentice Hall, Inc.

From Freeman and Herron (2007) "Evolutionary Analysis (4th edn)"

# How to analyze multiple sequences?

SCR-1476	TGAGACATAA	ACATTATATT	TTTCAATTAG	TCATTAAT	AAAAAATTCT
SCR-1640	TGAGACATAA	ACATTATATT	TTTCAATTAG	TCATTAAT	AAAAAATTCT
SCR-20	TGAGACATAA	ACATTATATT	TTTCAATTAG	TCATTAAT	AAAAAATTCT
SCR-6090	TGAGACATAA	ACATTATATT	TTTCAATTAG	TCATTAAT	AAAAAATTCT
SCR-6622	TGAGACATAA	ACATTATATT	TTTCAATTAG	TCATTAAT	AAAAAATTCT
SCR-6665	TGAGACATAA	ACATTATATT	TTTCAATTAG	TCATTAAT	AAAAAATTCT
SCR-6680	TGAGACATAA	ACATTATATT	TTTCAATTAG	TCATTAAT	AAAAAATTCT
SCR-6780	TGAGACATAA	ACATTATATT	TTTCAATTAG	TCATTAAT	AAAAAATTCT
SCR-6807	TGAGACATAA	ACATTATATA	TTTCAATTAG	TCATTAAT	AAAAAATTCT
SCR-6876	TGAGACATAA	ACATTATATT	TTTCAATTAG	TCATTAAT	AAAAAATTCT
SCR-6918	TGAGACATAA	ACATTATATT	TTTCAATTAG	TCATTAAT	AAAAAATTCT
SCR-915	TGAGACATAA	ACATTATATT	TTTCAATTAG	TCATTAAT	AAAAAATTCT
SCR-917	TGAGACATAA	ACATTATATT	TTTCAATTAG	TCATTAAT	AAAAAATTCT

# Contents: molecular population genetics

0. Genome-wide polymorphism and next-generation sequencers
1. nucleotide diversity  
selective sweep / hitchhiking, and balancing selection
2. Tajima's  $D$
3. positive selection on replacement mutations

# Nucleotide diversity $\pi$ (Nukleotid Diversität) $(\theta_\pi, \theta_T)$

Basic index of molecular population genetics

Average proportion of pairwise differences between the sequences  
 (Durchschnittliches Verhältnis von paarweisen Unterschieden zwischen den Sequenzen)

	* * * ..	(1)    (2)    (3)
(1) AGGCTGCATC	AGGCTGCATC	(1) ---
(2) .A.....	= AAGCTGCATC	(2) ---    ---
(3) .....T.C.	AGGCTGTACC	(3) ---    ---    ---

$$\pi = \sum_{i < j} \pi_{ij} / n_c = ?$$

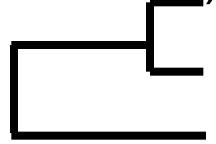
$$n_c = \frac{n(n - 1)}{2}$$

number of pairwise comparisons  
 (Anzahl von paarweisen Vergleichen)

$\pi_{ij}$  proportion of differences between i-th and j-th  
 (Verhältnis der Unterschiede zwischen i und j)

$$\pi = \sum_{i < j} \pi_{ij} / n_c = \frac{\pi_{12} + \pi_{13} + \pi_{23}}{3} = ?$$

Genealogy  
 (Stammbaum)



# nucleotide diversity $\pi$

Your own calculation would hopefully help intuitive understanding  
(on average, 2 out of 10 are different)

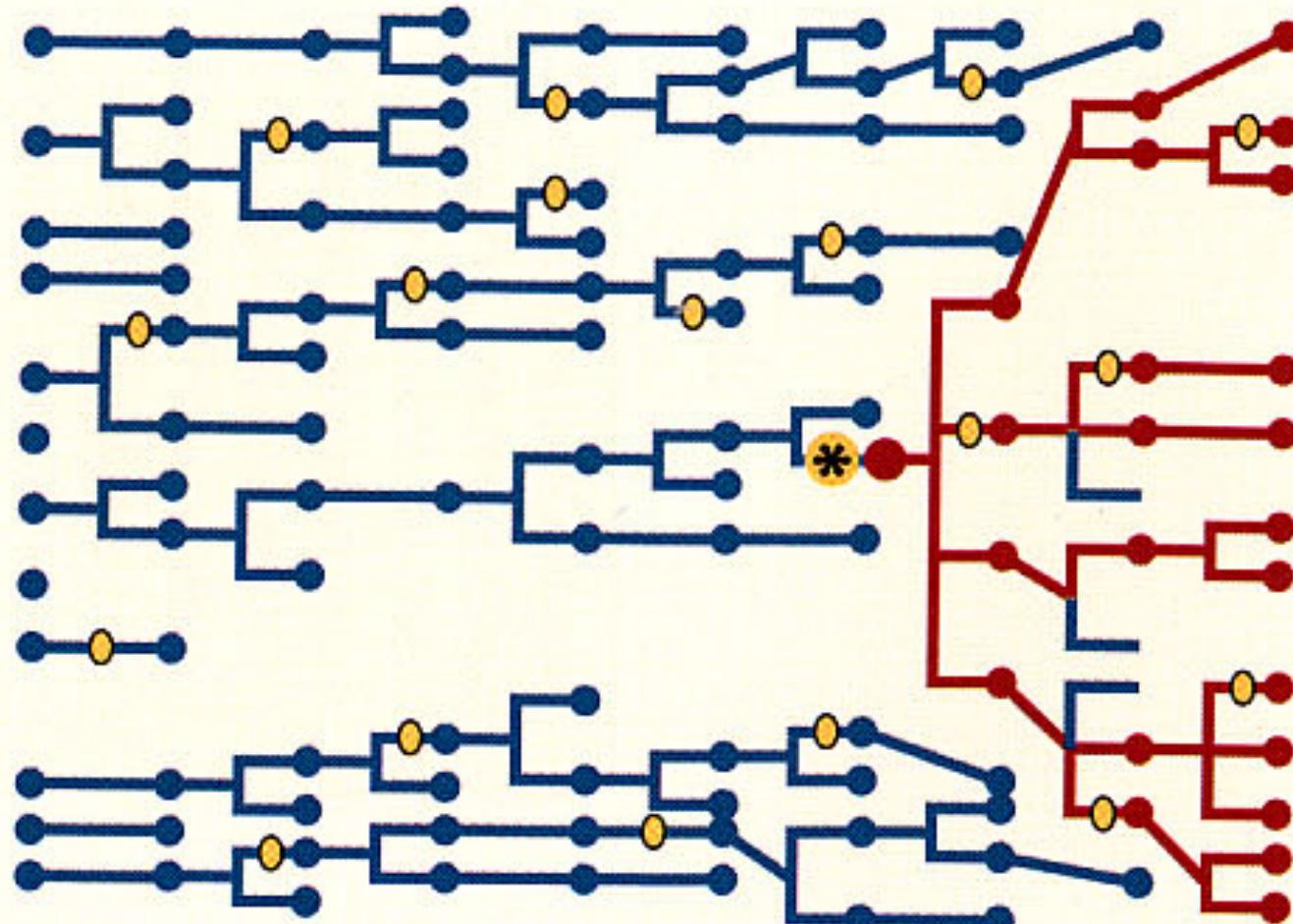
Genomic average of nucleotide diversity  $\pi$

*Homo sapiens*: ~0.0008

*Arabidopsis thaliana* and *Drosophila melanogaster*: ~0.007

# Coalescent theory: coalescence and mutations (Koaleszenz bzw. Vereinigungs- Theorie – Koaleszenz und Mutationen)

Futuyma, Evolution, Sinauer 2005 p. 289



AGGCTGCATC  
.....T...  
.....  
.A.T....  
.A.....  
.....  
.....C...  
.....  
.....T...  
.....T...

Natural selection can be classified into 2 types based on the pattern of polymorphism

## **1. positive selection**

(directional selection, Darwinian selection):

selection for an allele that increases fitness.

Positive selection results in selective sweep.

## **2. balancing selection**

a form of natural selection that maintains polymorphism at a locus within a population.

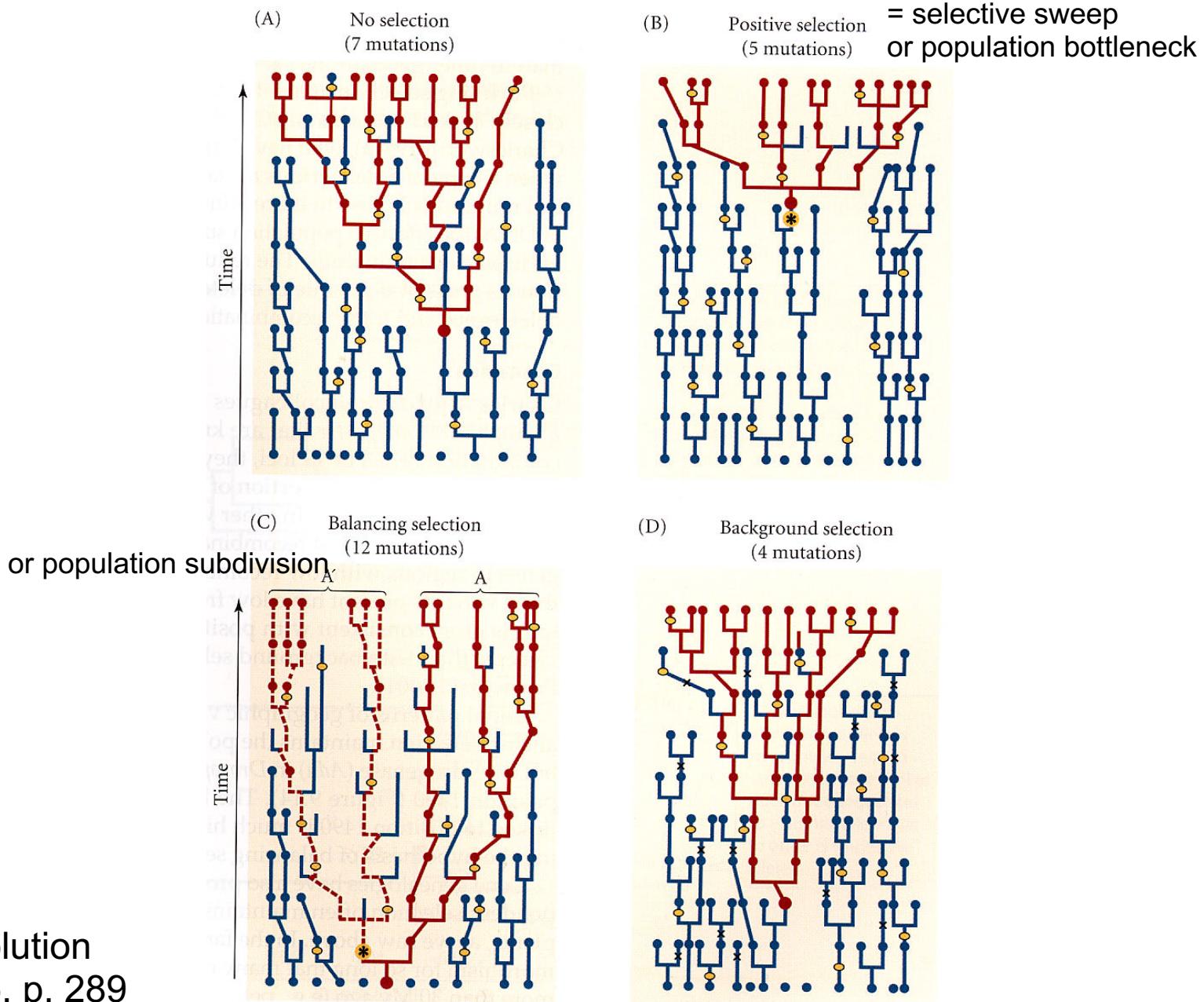
This includes heterozygosity advantage, negative frequency dependent selection, local adaptation, etc. Diversifying selection can be used in a similar meaning.

cf. purifying selection (negative selection):

elimination of deleterious alleles from a population

le gene is eliminated by a selective sweep neutral mutations occur among

- Neutral mutation
- ✳ Advantageous mutation
- ✗ Deleterious mutation

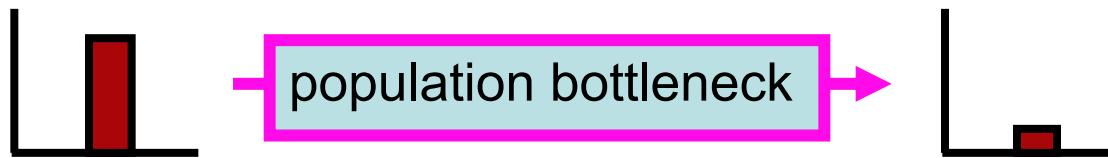


# Natural selection vs. population processes (or demography) Natürliche Selektion vs. Demografie

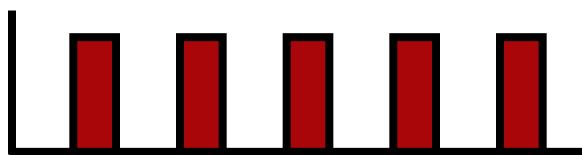
Gene  
genealogy



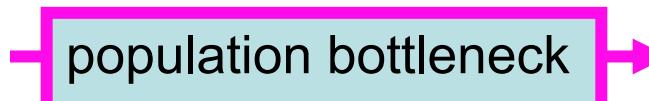
nucleotide  
diversity  $\pi$



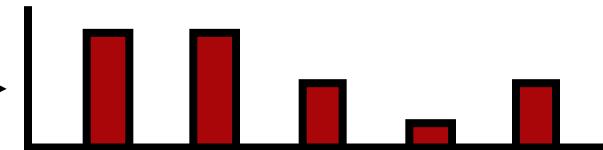
## Power of Genomics



chromosome



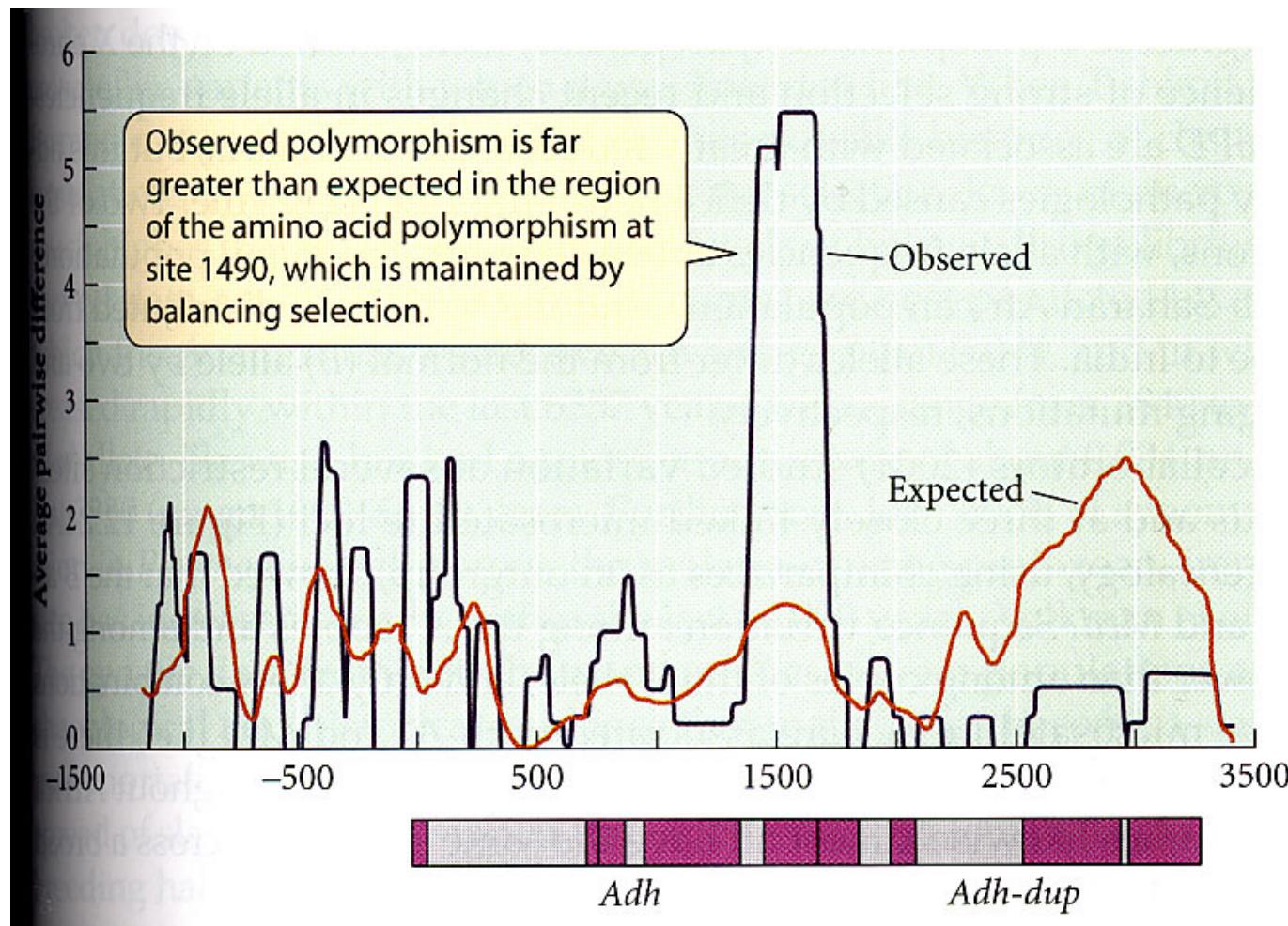
## recombination



$\pi$  lower than genomic average and neighbors  
→ positive selection

# balancing selection

(“Selektion die Polymorphismen erhält”)



# Contents: molecular population genetics

0. Genome-wide polymorphism and next-generation sequencers
1. nucleotide diversity  
selective sweep / hitchhiking, and balancing selection
2. Tajima's  $D$
3. positive selection on replacement mutations

# $\theta$ , another index ( $\theta$ , ein weiterer Index)

(nucleotide polymorphism based on polymorphic site /  
Nukleotid Polymorphismus basierend auf der Anzahl  
voneinander verschiedener Nukleotide an einer bestimmten  
Stelle in einem Sequenzvergleich)

	*      *    *	(1)	(2)	(3)
(1) AGGCTGCATC	AGGCTGCATC	(1)	---	---
(2) .A.....	= AAGCTGCATC	(2)	---	---
(3) .....T.C.	AGGCTGTACC	(3)	---	---

$$\pi = \sum_{i < j} \pi_{ij} / n_c = ?$$

$$\theta = s / \sum_{k=1}^{n-1} \frac{1}{k} = ?$$

s: the proportion of polymorphic sites (or nucleotide polymorphism) observed in the sample (Das Verhältnis von "voneinander verschiedenen Stellen" bzw. Nukleotid Polymorphismen, die in einem Sequenzvergleich beobachtet werden können.)

# Test of neutrality: Tajima's $D$

## (Testen auf Neutralität: Tajima's $D$ )

example: balancing selection

Tajima Genetics 123, p. 229, 1989

(1) AGGCTGCATC

(2) .....

(3) .A.....C.

(4) .A.....C.

$$\pi = \sum_{i < j} \pi_{ij} / n_c = ?$$

$$\theta = s / \sum_{k=1}^{n-1} \frac{1}{k} = ?$$

$$D = \frac{\pi - \theta}{\text{standard\_deviation\_of\_}(\pi - \theta)}$$

positive or negative?

# Test of neutrality: Tajima's $D$

example: positive selection

(1) AGGCTGCATC

(2) .....

(3) .....

(4) .A.....C.

$$\pi = \sum_{i < j} \pi_{ij} / n_c = ?$$

$$\theta = s / \sum_{k=1}^{n-1} \frac{1}{k} = ?$$

$$D = \frac{\pi - \theta}{\text{standard\_deviation\_of\_}(\pi - \theta)}$$

positive or negative?

# $\theta$ is less affected by frequency

example: balancing selection

- (1) AGGCTGCATC
- (2) .....
- (3) .A.....C.
- (4) .A.....C.

example: positive selection

- (1) AGGCTGCATC
- (2) .....
- (3) .....
- (4) .A.....C.

$$\theta = s / \sum_{k=1}^{n-1} \frac{1}{k} =$$

$$D = \frac{\pi - \theta}{\text{standard\_deviation\_of } (\pi - \theta)}$$

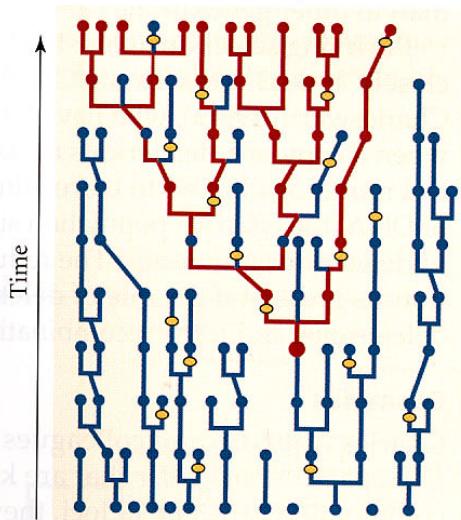
Less than zero when  $\pi$   
is low due to positive selection

Intuitively: singleton vs. shared mutation  
DnaSP and other free softwares

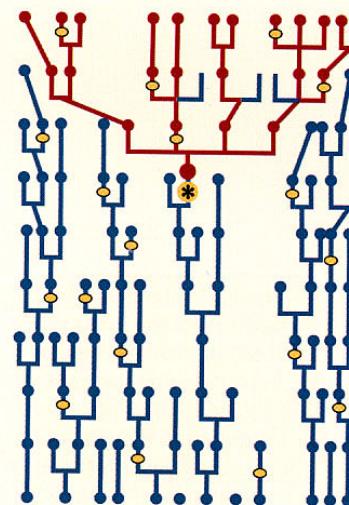
the gene is eliminated by a selective sweep neutral mutations occur among

- Neutral mutation
- \* Advantageous mutation
- x Deleterious mutation

(A) No selection  
(7 mutations)



(B) Positive selection  
(5 mutations)



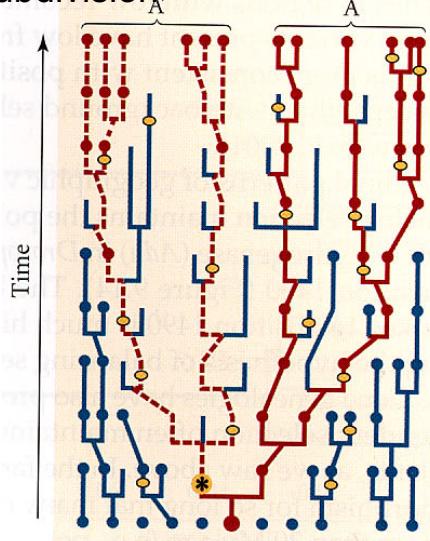
= selective sweep  
or population expansion

$D < 0$   
more  
singletons

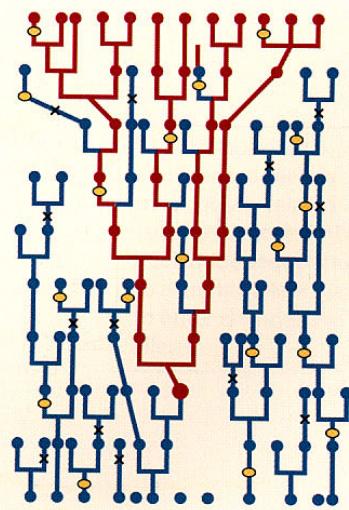
or population subdivision

$D > 0$   
fewer  
singletons

(C) Balancing selection  
(12 mutations)



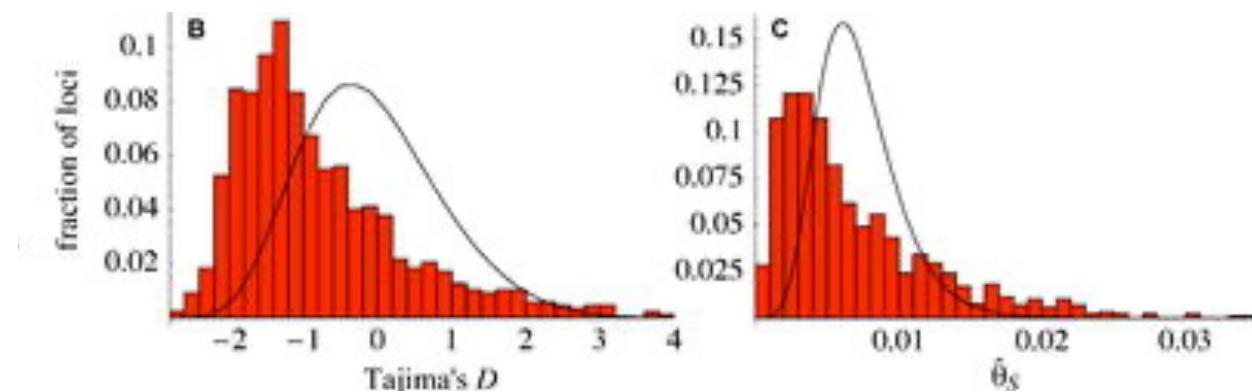
(D) Background selection  
(4 mutations)



# empirical distribution: average may not be zero

876 short fragments in a sample of 96 individuals of *Arabidopsis thaliana*

Nordborg et al. The pattern of polymorphism in *Arabidopsis thaliana*. PLoS Biol. 2005 Jul;3(7):e196



# Contents: molecular population genetics

0. Genome-wide polymorphism and next-generation sequencers
1. nucleotide diversity  
selective sweep / hitchhiking, and balancing selection
2. Tajima's  $D$
3. positive selection on replacement mutations

# Synonymous (silent-site) and nonsynonymous (replacement) mutations

First base	Second base		Third base
	U	C	
U	UUU Phenylalanine	UCU Serine	U
	UUC Phenylalanine	UCC Serine	C
	UUA Leucine	UCA Serine	A
	UUG Leucine	UCG Serine	G

Figure 7-21a Evolutionary Analysis, 4/e  
© 2007 Pearson Prentice Hall, Inc.

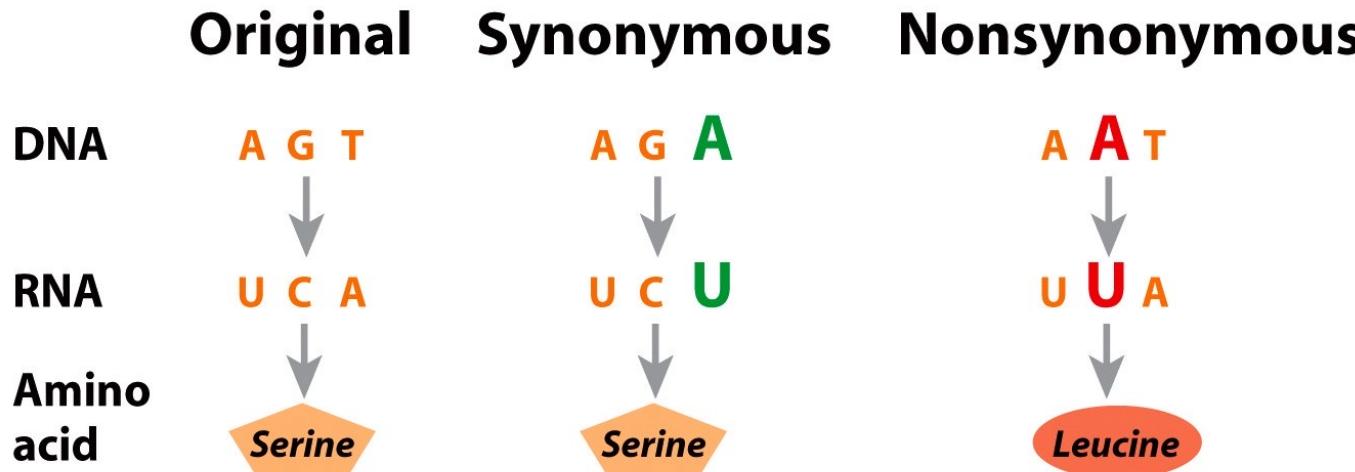


Figure 7-21b Evolutionary Analysis, 4/e  
© 2007 Pearson Prentice Hall, Inc.

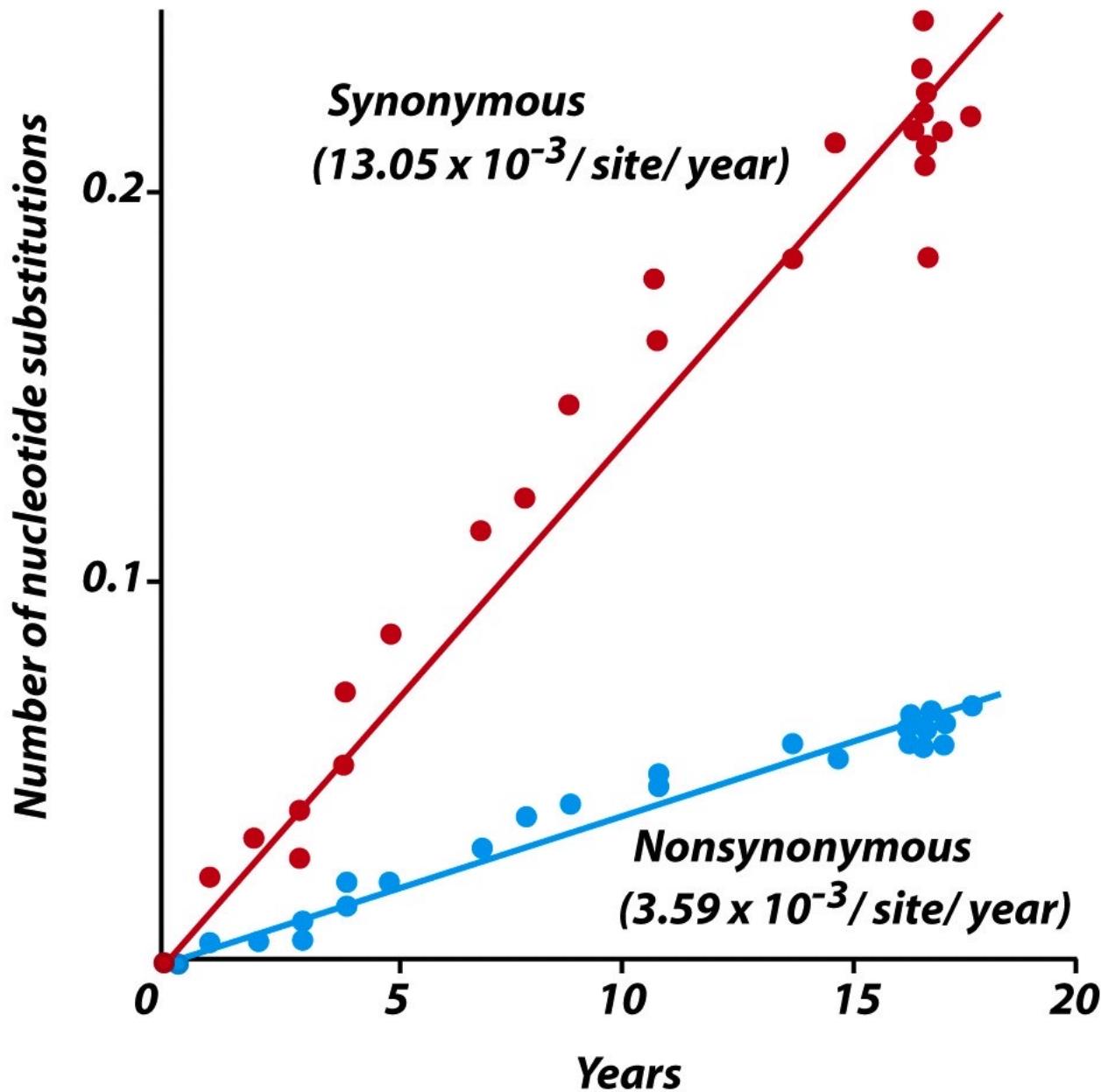


Figure 7-21c Evolutionary Analysis, 4/e

© 2007 Pearson Prentice Hall, Inc.

$$\omega = K_a/K_s = d_N / d_S$$

$d_N / d_S < 1$  when replacements are deleterious

$d_N / d_S = 1$  when replacements are neutral

$d_N / d_S > 1$  when replacements are advantageous  
(positive selection)

NOTE:  $d_N / d_S > 1$  indicates a strong evidence for natural selection,  
but its sensitivity is low

Sharp (1997) points out that “it is extremely conservative.”

Only when many amino acids changed, it can be more than 1.

$d_N / d_S$  is often used  
if only one sequence from each species are available

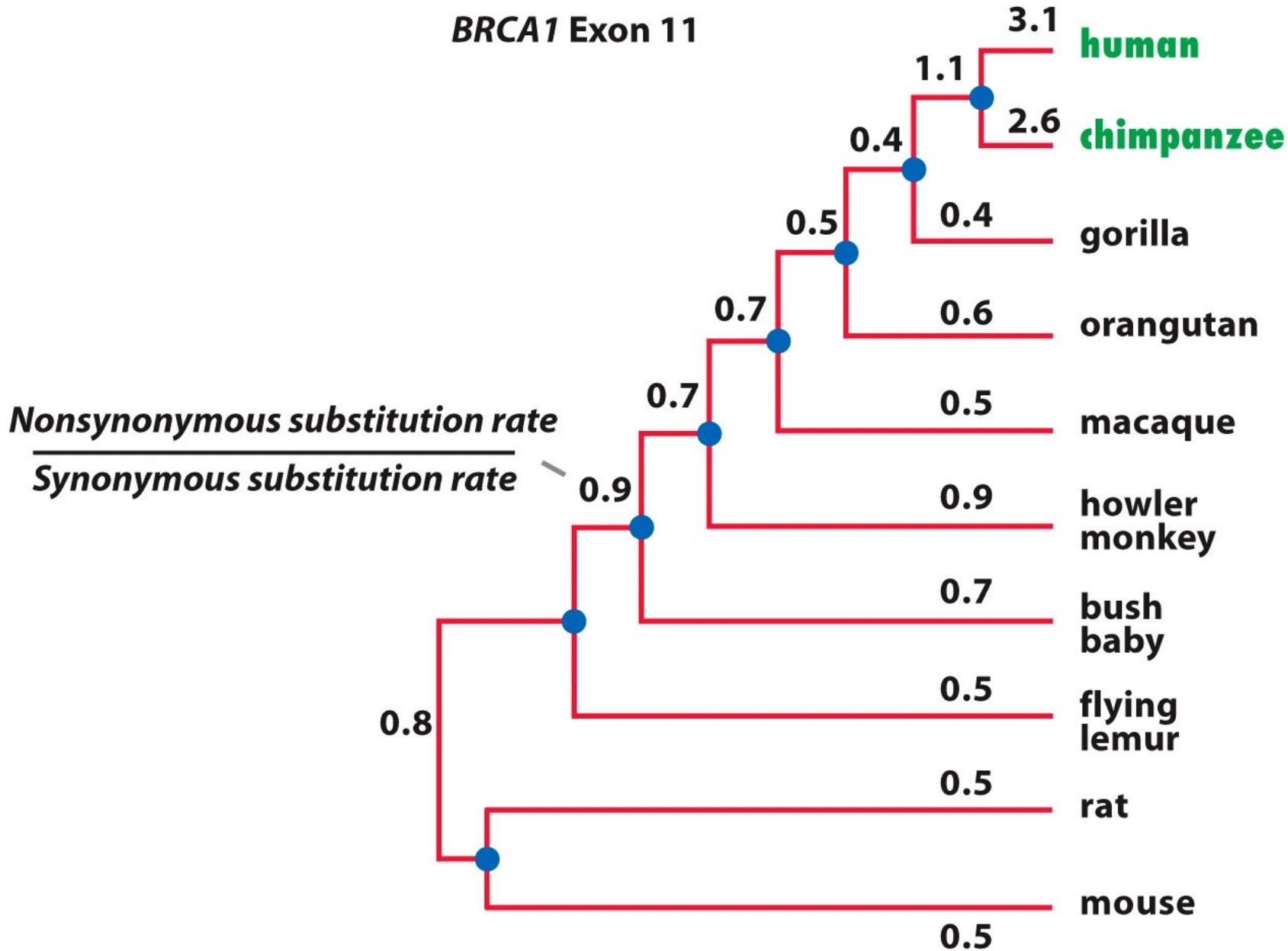


Figure 7-23 Evolutionary Analysis, 4/e  
 © 2007 Pearson Prentice Hall, Inc.

**Table 7.2 Studies that confirm positive selection on replacement mutations**

Although this table lists just a few examples, it underscores a general point: Evidence for positive selection is particularly strong in genes that code for proteins involved in disease resistance, reproductive conflict, interactions between symbionts, and the development of novel traits.

Gene	Species	Rationale	Reference
MHC Class II	Humans	Strong selection for divergence among antigen-recognition proteins	Hughes, A.L., and M. Nei. 1989. <i>Proceedings of the National Academy of Sciences USA</i> 86: 958–962.
Semenogelin II gene	Primates	Strong selection for semen coagulation in species with intense sperm competition	Dorus, S. D., et al. 2004. <i>Nature Genetics</i> 12: 1326-1329.
Lysin protein and receptor	Abalone	Strong selection on species-specific egg-recognition proteins on sperm	Swanson, W.J., and V.D. Vacquier. 1998. <i>Science</i> 281: 710–712.
Eosinophil cationic protein	Primates	Strong selection on a recently duplicated gene involved in disease resistance	Zhang, J., H.F. Rosenberg, and M. Nei, 1998. <i>Proceedings of the National Academy of Sciences USA</i> 95: 3708–3713.
Self-incompatibility loci	Tomato family plants	Strong selection for divergence among proteins involved in self-fertilization	Clark, A.G., and T.-H. Kao. 1991. <i>Proceedings of the National Academy of Sciences USA</i> 88: 9823–9827.
Multicolored fluorescent proteins	Reef-building corals	Evolving algal symbionts impose selection on ability of host corals to regulate them.	Field, S. F., et al. 2006. <i>Journal of Molecular Evolution</i> 62: 332-339.
ASPM	Humans and great apes	Selection for increased brain size.	Evans, P. D., et al. 2004. <i>Human Molecular Genetics</i> 13: 489-494.

# Case study: Evolution of language

Was there positive selection to have language?

How many genes were necessary to obtain language?

How and when did the language evolve?

Did Neanderthal have language,  
which diverged 3-400,000 years ago?

Limited evidence: archaeology, fossil of skeletons

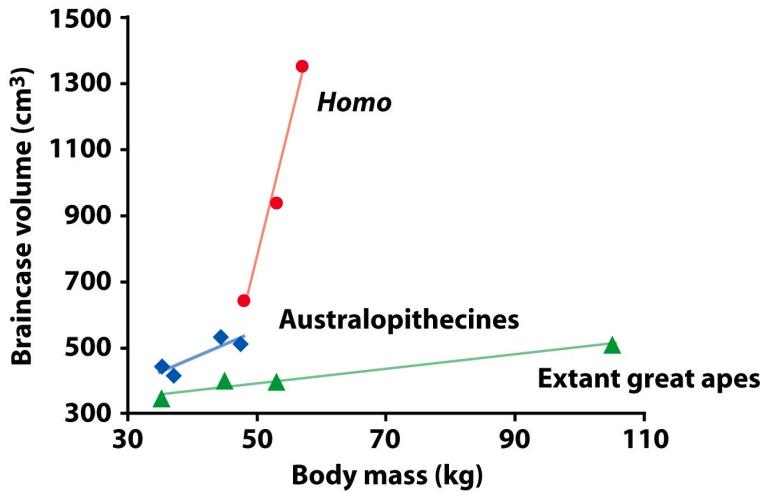


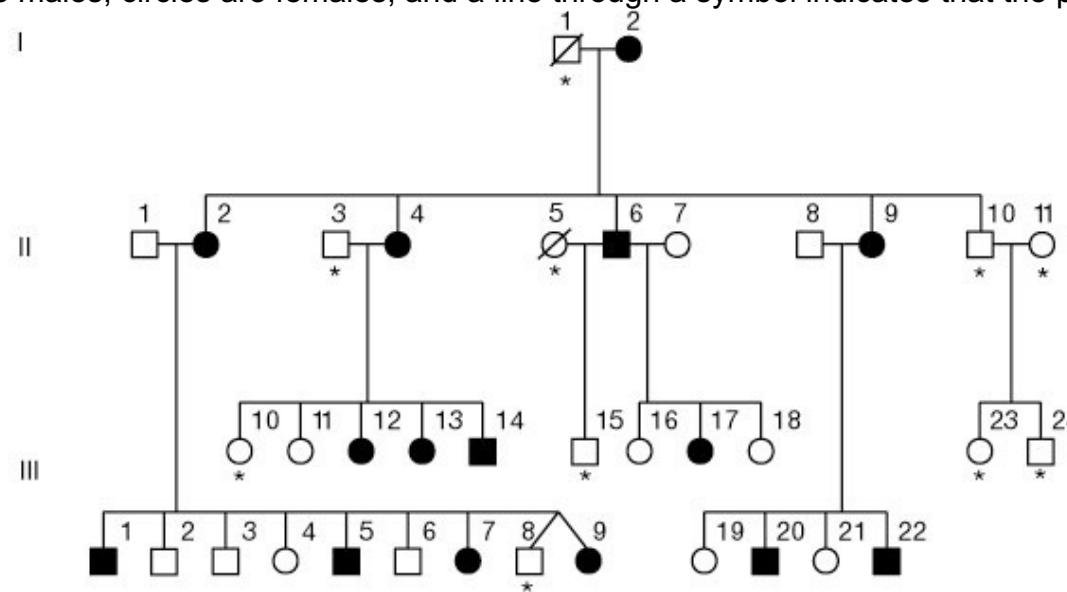
Figure 20-31 Evolutionary Analysis, 4/e  
© 2007 Pearson Prentice Hall, Inc.

Figure 20-32 Evolutionary Analysis, 4/e  
© 2007 Pearson Prentice Hall, Inc.

# Candidate gene from human disease genes: FoxP2, necessary for the development of Broca's area and language

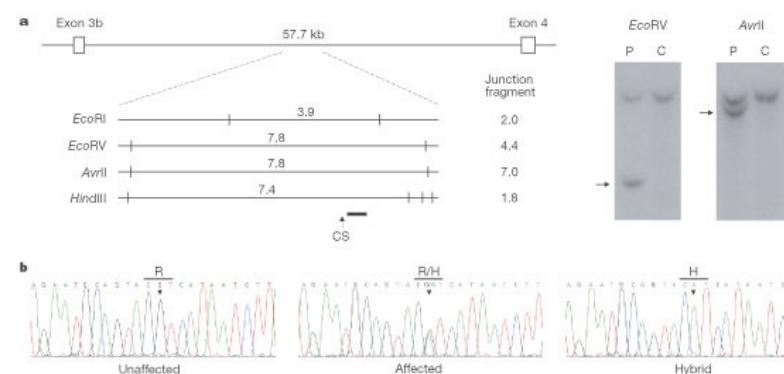
**FIGURE 1. Pedigree of the KE family.**

Affected individuals are indicated by filled symbols. Asterisks indicate those individuals who were unavailable for genetic analyses. Squares are males, circles are females, and a line through a symbol indicates that the person is deceased.



## A forkhead-domain gene is mutated in a severe speech and language disorder

Lai et al. *Nature* 413, 519-523, 2001



# The signatures of recent positive selection

## (Anzeichen von rezenter positiver Selektion)

Table 1 Variation at the *FOXP2* locus in humans

No. of chromosomes sequenced	40
Length covered (double stranded, all individuals)	14,063 bp
Divergence from the chimp sequence*	0.87%
No. of variable positions	47
Singletons (no. of variable sites occurring at frequency 1 and 39)	31
$\theta_W$ (nucleotide diversity based on the no. of polymorphic sites)	0.079%
$\theta_\pi$ (mean nucleotide diversity)	0.03%
$\theta_H$ (nucleotide diversity with more weight given to alleles at high frequency <sup>17</sup> )	0.117%
$D$ ( $P < 0.01$ )†	-2.20
$H$ ( $P < 0.05$ )‡	-12.24

$\pi \rightarrow$

\*The corresponding value for the orang-utan is 2.5.

†A negative  $D$  value indicates a relative excess of low-frequency alleles<sup>15</sup>.

‡A negative  $H$  value indicates a relative excess of high-frequency derived alleles<sup>17</sup>.

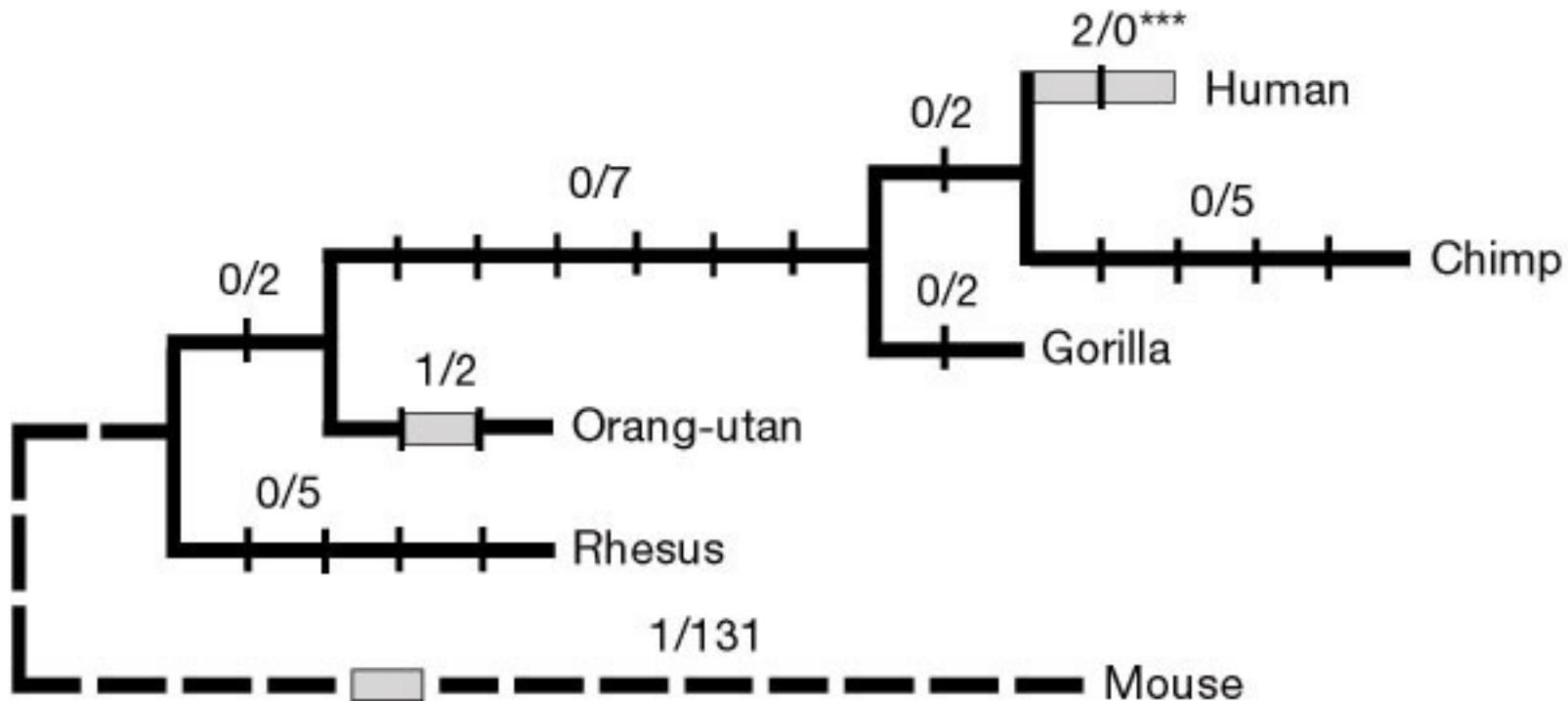
### Molecular evolution of *FOXP2*, a gene involved in speech and language

Wolfgang Enard, Molly Przeworski, Simon E. Fisher, Cecilia S. L. Lai, Victor Wiebe, Takashi Kitano, Anthony P. Monaco and Svante Pääbo

*Nature* **418**, 869-872, 2002

Zheng et al., *Genetics* 162, p. 1825 (2002)

high Ka/Ks in the human lineage  
(hoher Ka/Ks in der menschlichen Abstammung)



# When did the language appear?

Krause et al. (2007) Curr. Biol. 17: 1908

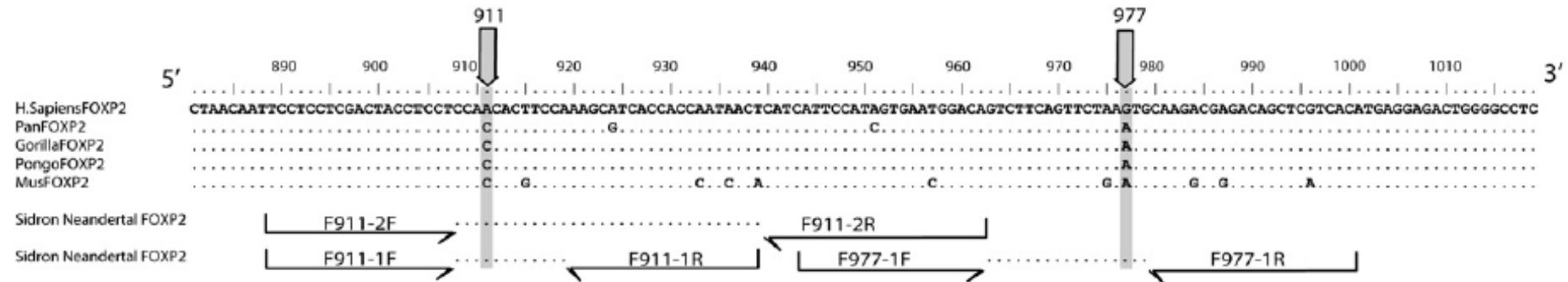


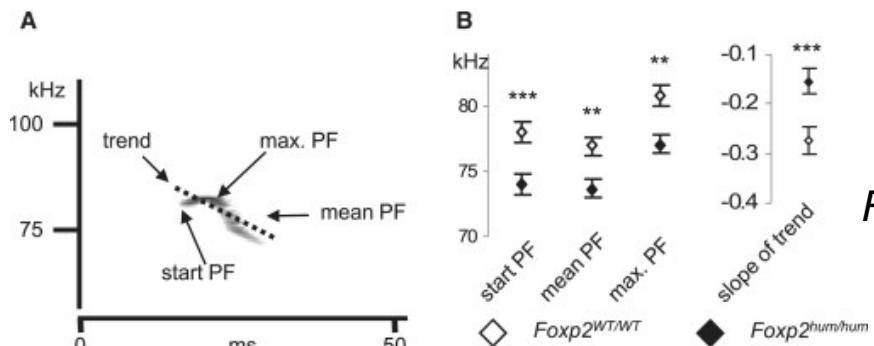
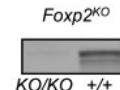
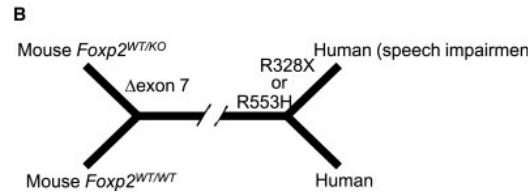
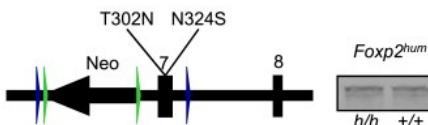
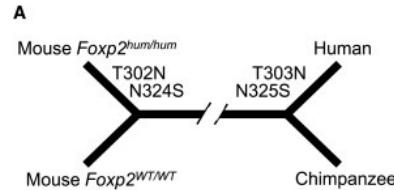
Figure 1. Sequence Alignment of Nucleotide Positions 880–1020 from the *FOXP2* Gene

The two nonsynonymous nucleotide substitutions on the human lineage are indicated by arrows. Identical positions in the alignment are given as dots. The three primer pairs used to retrieve the two substitutions from the El Sidrón Neandertals are indicated by arrows.

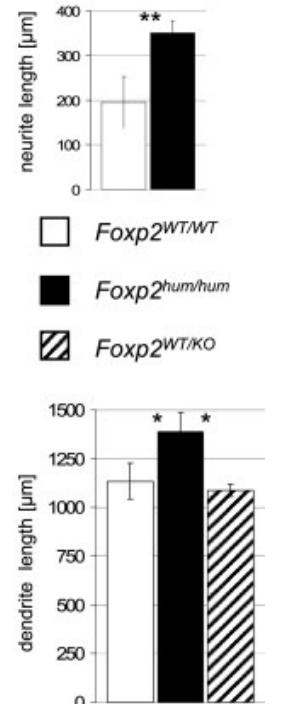
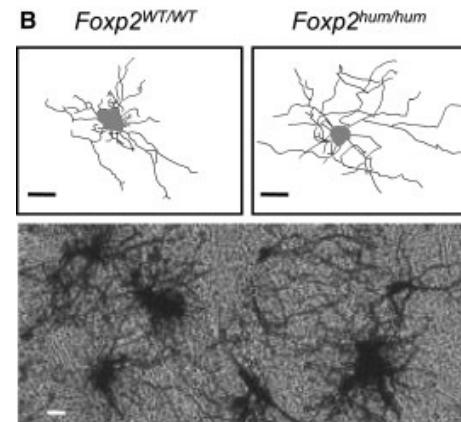
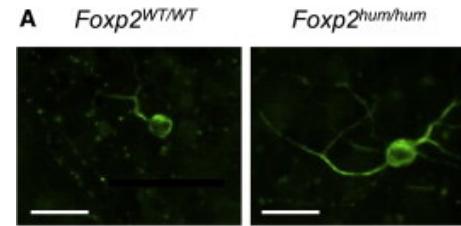
Time estimate based on coalescent simulation:  
with approximate 95% confidence intervals of 0 and 120,000 years ago  
---> Time estimate has a number of assumptions and may not be accurate

Note: There is little evidence that this gene really contributed to the evolution of language. Even if so, multiple genes must have been necessary.

# Transgenic mouse with human FoxP2

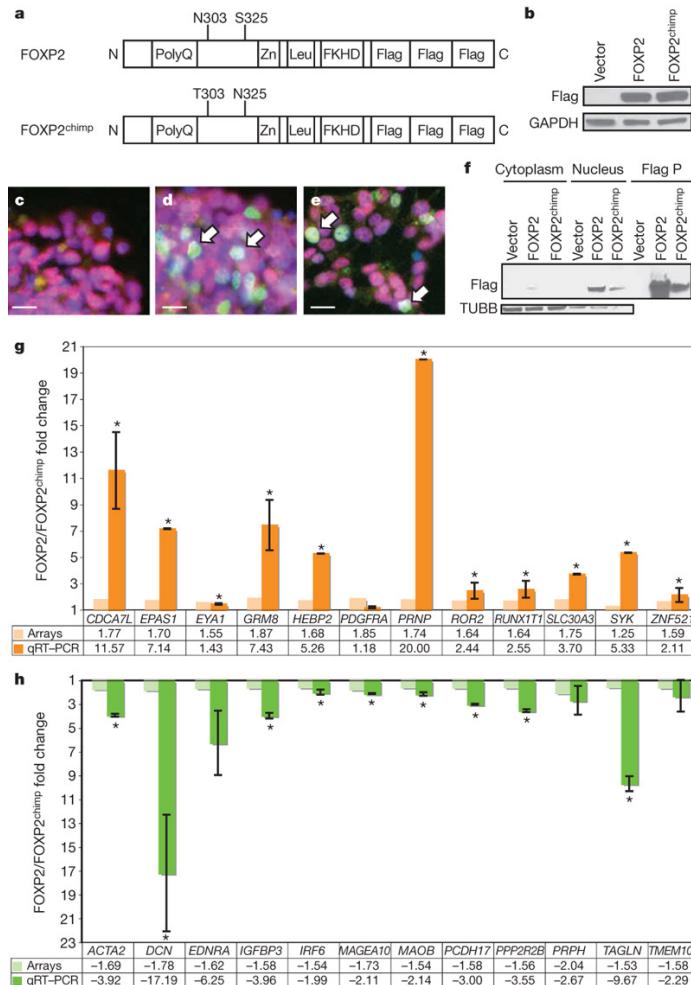


$Foxp2^{hum}$  Affects the Structure of Pup Isolation Calls



$Foxp2^{hum}$  Increases the Length of Dendritic Trees

# FOXP2 and FOXP2<sup>chimp</sup> differentially regulate genes in neural cells.



G Konopka et al. *Nature* 462, 213-217 (2009)  
doi:10.1038/nature08549

# Summary

- What was revealed about evolution by using molecular data?
- Why are model species studied in addition to *Homo sapiens*?

# Exercise example: polyploid species

- Genome duplication (polyploidization) is prevalent in animals, fungi and plants (>35%)



Ohno 1970, Swalla, *Heredity* 2006



Leitch & Leitch *Science* 2008

# A new species during the past 150 years in Switzerland

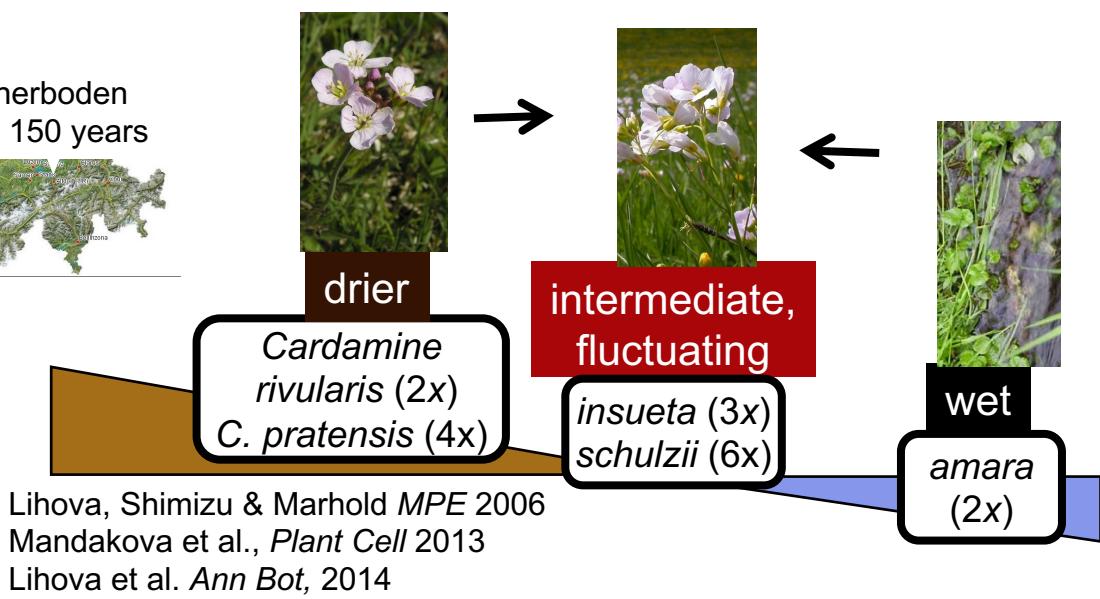
A textbook example in Switzerland:



Urnerboden  
past 150 years



(picture by National Geographic)



drier  
*Cardamine rivularis* (2x)  
*C. pratensis* (4x)

intermediate,  
fluctuating  
*insueta* (3x)  
*schulzii* (6x)

wet  
*amara*  
(2x)

# A new species during the past 150 years in Switzerland



# Challenge in the polyploid transcriptome: quantifying ratio of homologous pairs

Major obstacle: difficult to separate homeologous pairs due to high homology

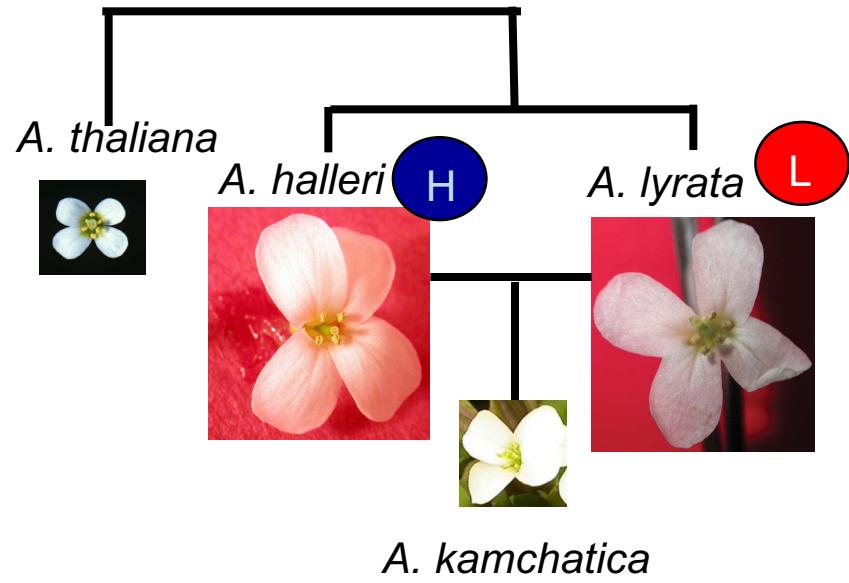
1. Genome assembly of polyploid is difficult
  - > new sequencers are solving it
2. Transcriptome
  - Designing homeolog-specific primer is time-consuming
  - Few tools to deal with genome-wide pattern

(SNP-based method: Shi et al. *Nature Comm* 2012, Page et al. *G3*, 2013)

-> New bioinformatic workflow HomeoRoq

# *Arabidopsis kamchatica*

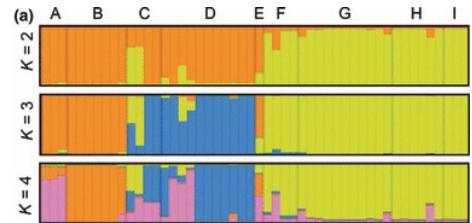
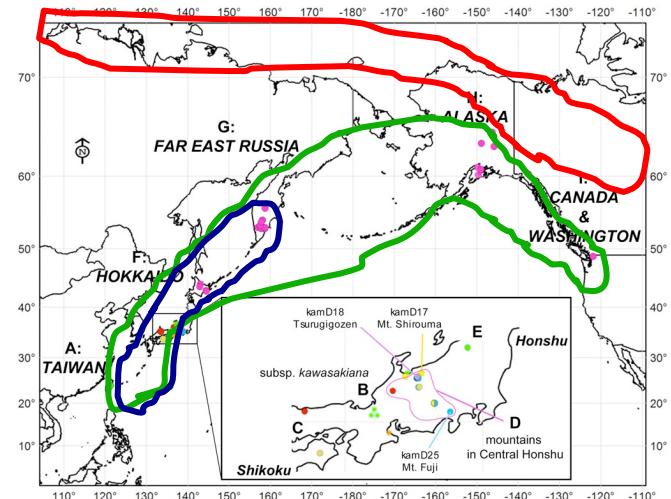
## a broadest ecological niche in the genus *Arabidopsis*



ssp. *kamchatica*, Perennial  
0-3000 m alt

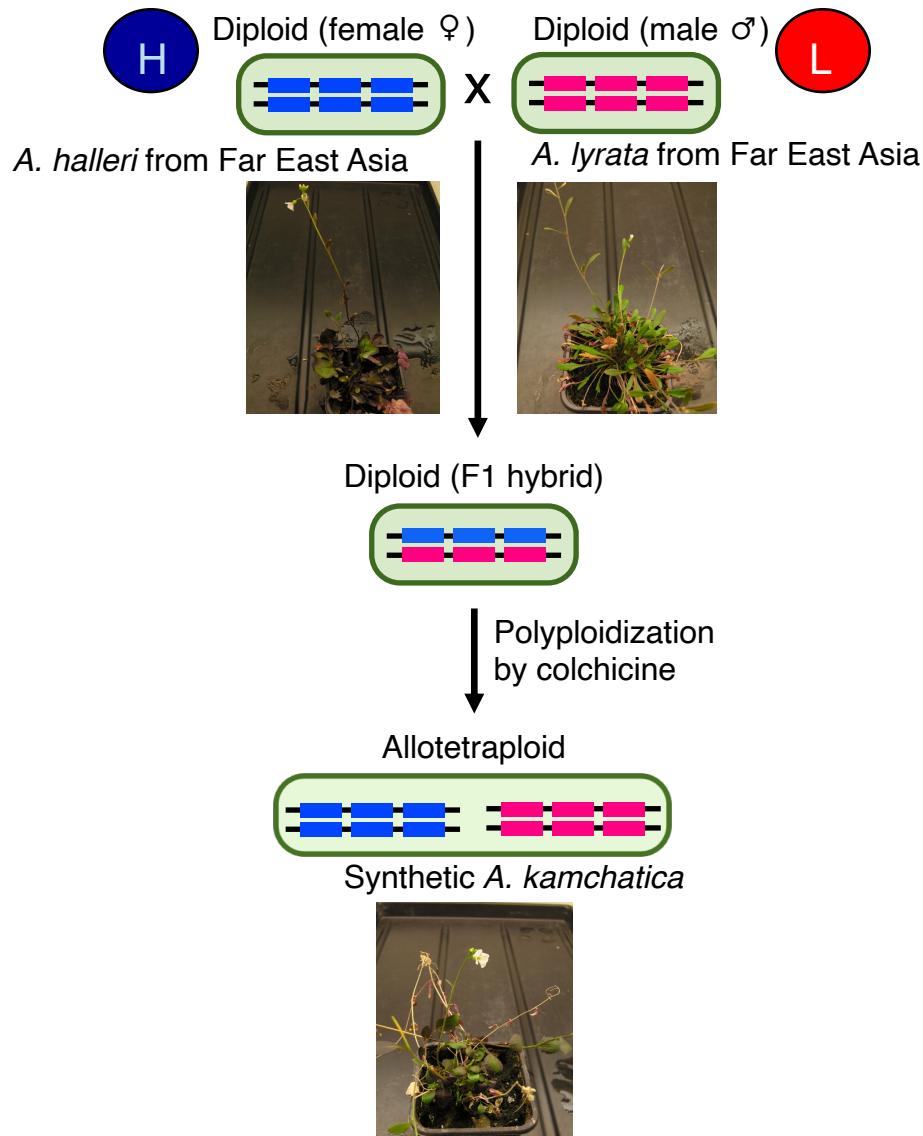


ssp. *kawasakiiana*  
Annual  
hot sea- or lake shore  
0-85 m alt in Japan

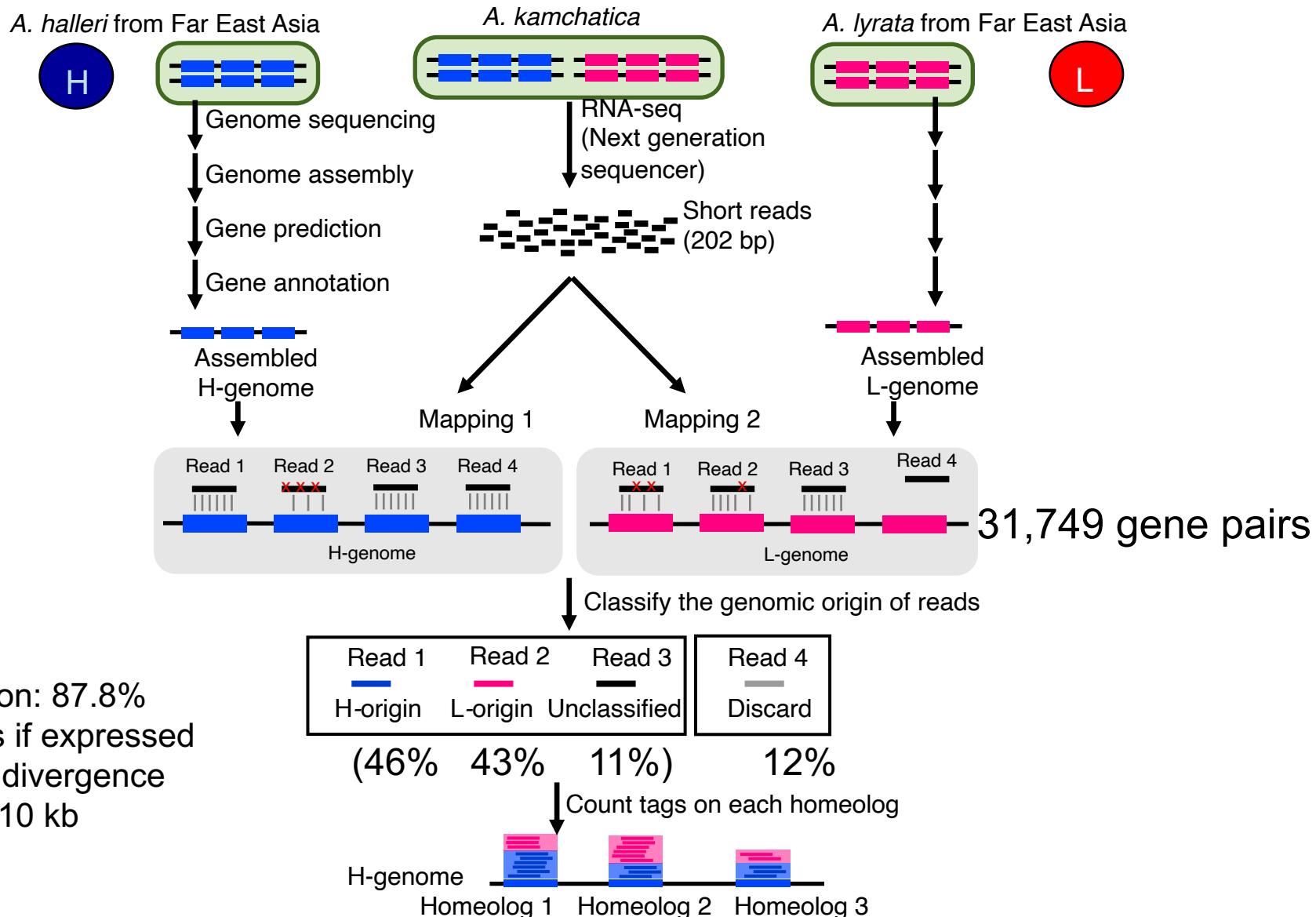


Shimizu-Inatsugi et al. *Mol Ecol* 2009

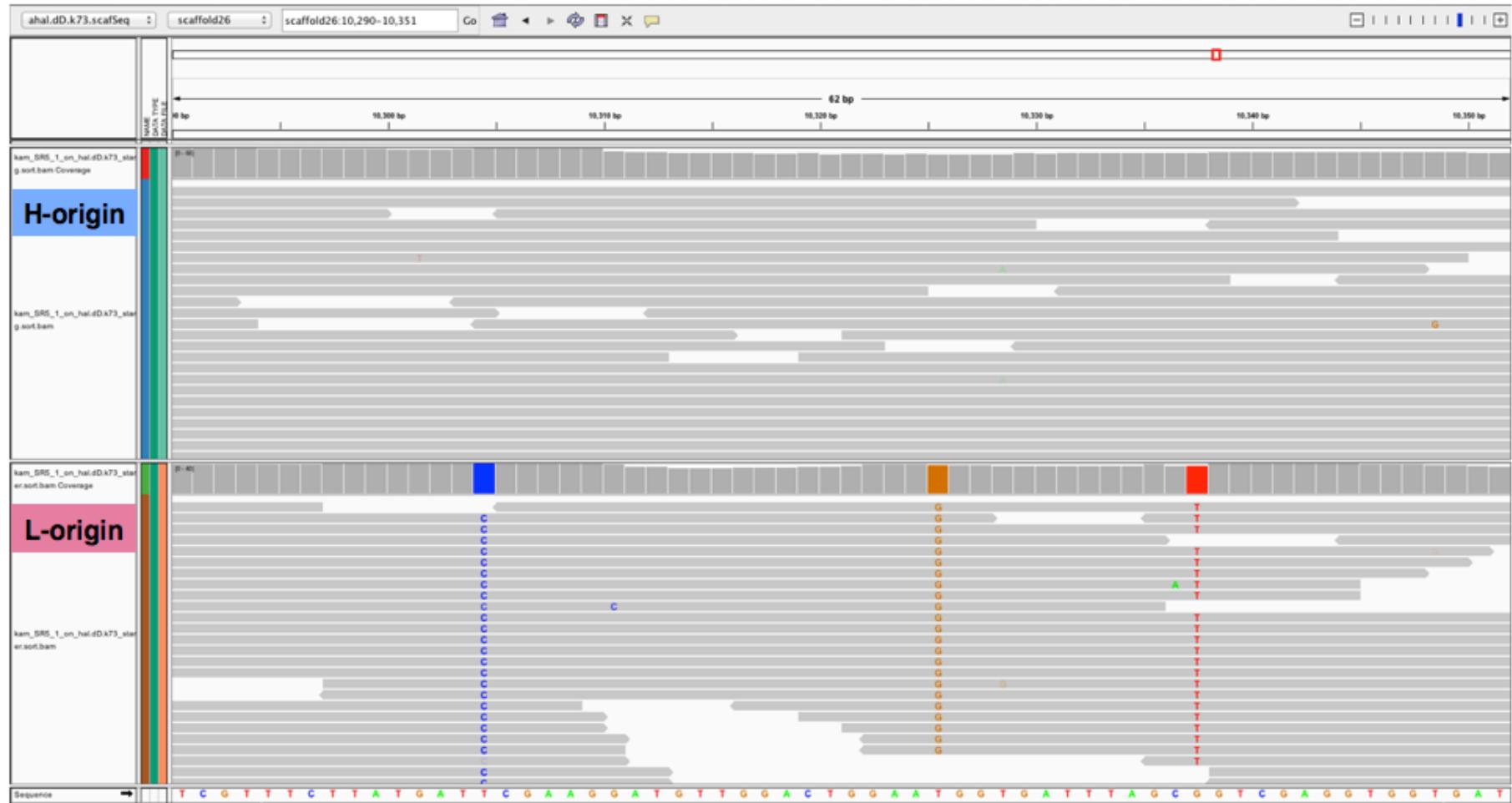
# Simplest polyploid model: Genome assembly of two parental species and synthetic allopolyploid



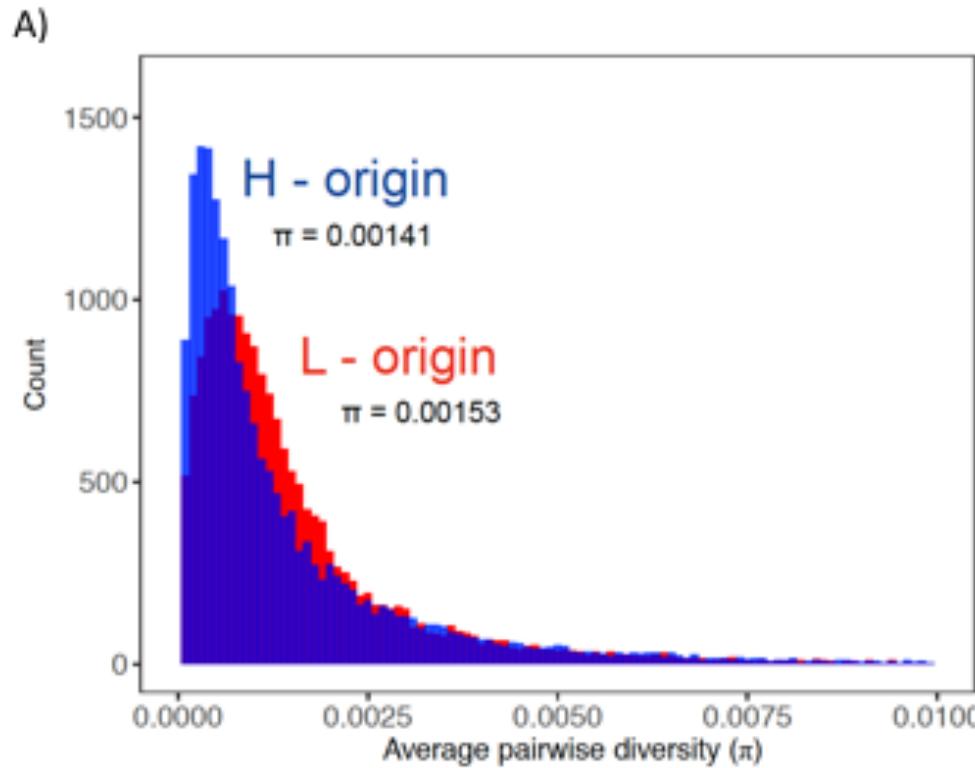
# HomeoRoq (Homeolog Ratio and Quantification): bioinformatic workflow using NGS



# Visualization: homeologs were well separated



# nucleotide diversity of polyploid species



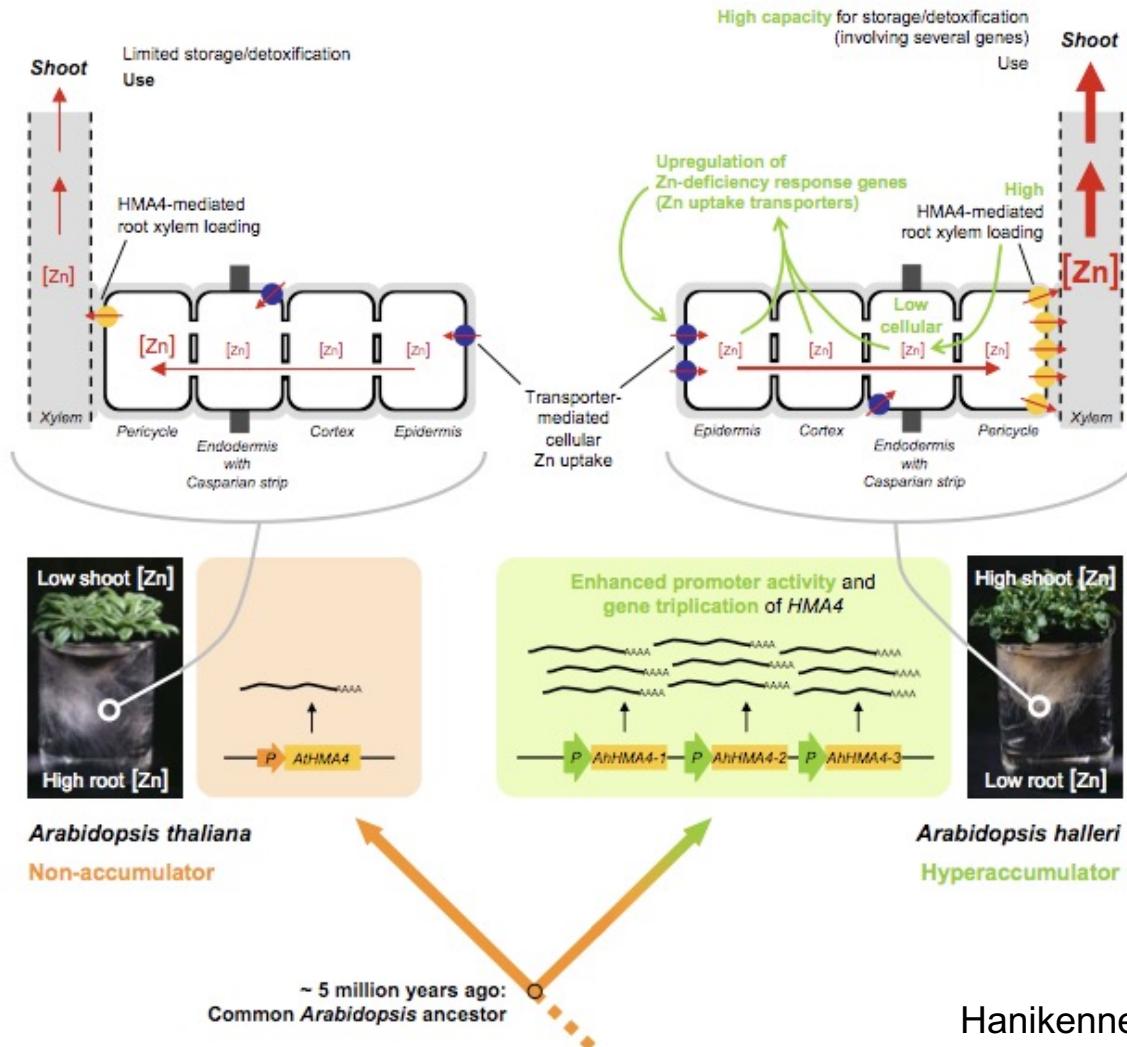
synonymous  $\pi$ : 0.0049 in H-subgenome, 0.0044 in L-subgenome

similar to *A. thaliana* (0.0059 - 0.007)

Several times lower than parental species (0.028 for *A. halleri* and *A. lyrata*)  
suggesting a bottleneck at the polyploidization events

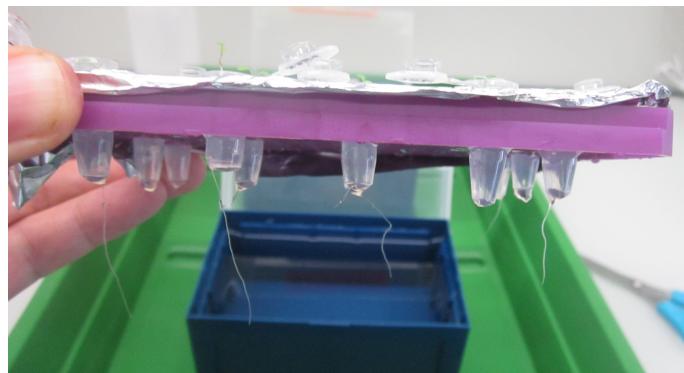
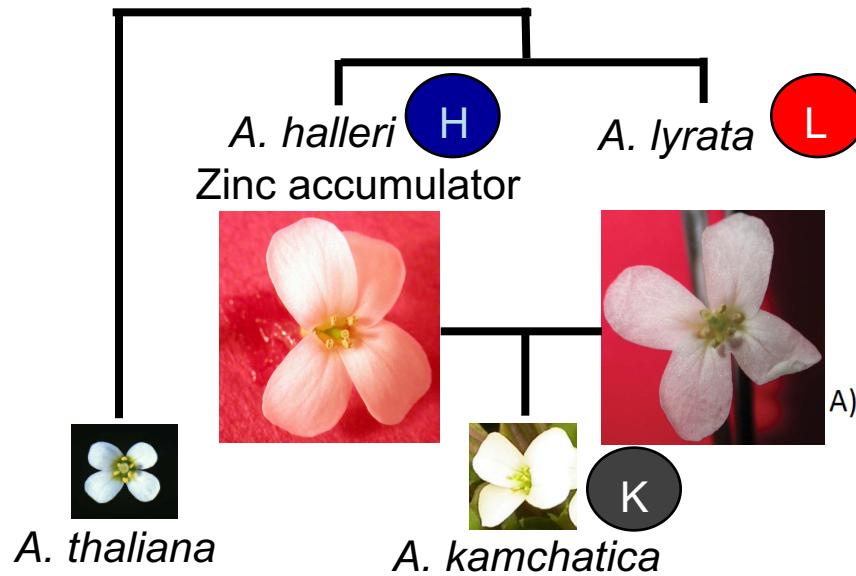
# Hyperaccumulation and phytoremediation

## *HMA4* transporter gene

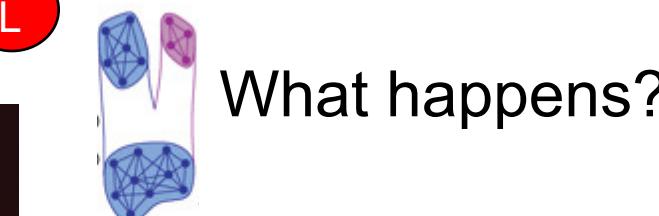


Hanikenne et al 2008 Nature

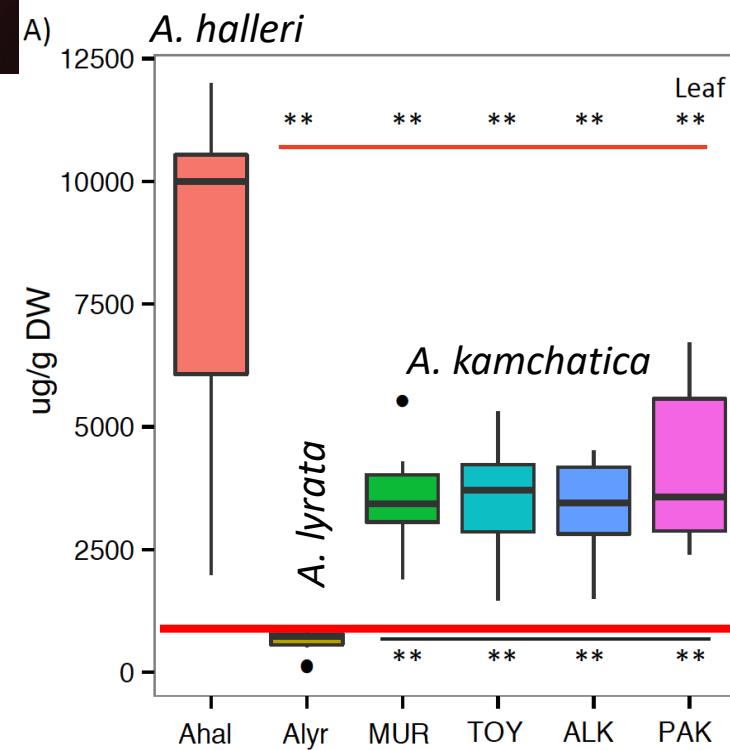
# Zinc hyperaccumulation to study both phenotype and gene expression level quantitatively



500  $\mu\text{M}$  Zinc for 7 days treatment



Paape et al.,  
*Mol Biol Evol* 2016



Insect  
defense

# Soil contamination



- Silver from Tada mine during the 16th century
- Major source to build the Osaka Castle when Hideyoshi Toyotomi unified Japan
- Closed in 1973



Zinc concentration of habitat soil ( $\text{mg kg}^{-1}$ )

*A. halleri*:

3.4-20 (nonmetalliferous) to 200-3,200 (metaliferous)

(Hanikenne et al. *PLoS Genet* 2013)

*A. kamchatica* 30 localities: 31-1,100

# Natural selection vs. population processes (or demography) Natürliche Selektion vs. Demografie

Gene

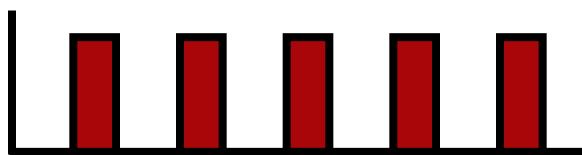
genealogy



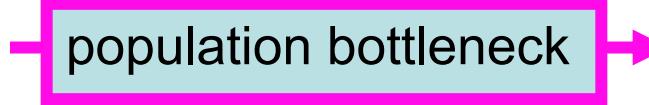
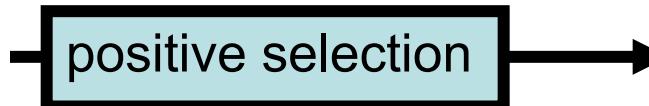
nucleotide  
diversity  $\pi$



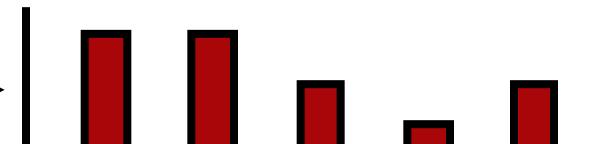
## Power of Genomics



chromosome



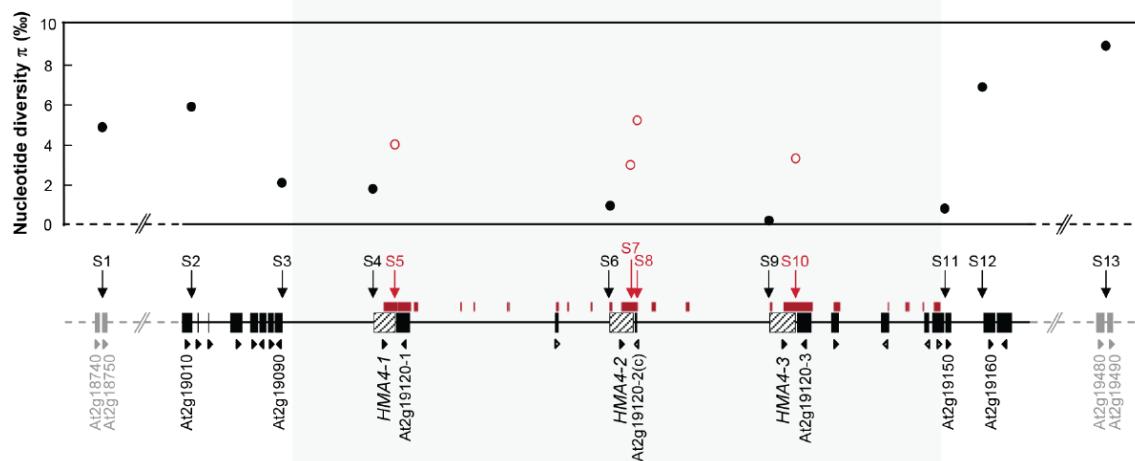
recombination



$\pi$  lower than genomic average and neighbors  
→ positive selection

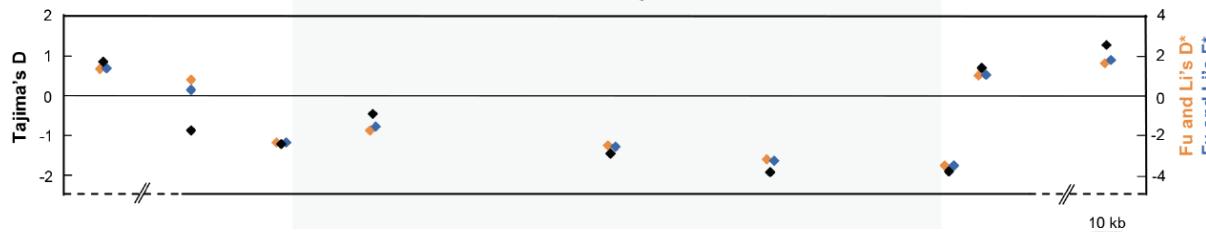
# HMA4 diversity statistics in the diploid parent *A. halleri*

**A**

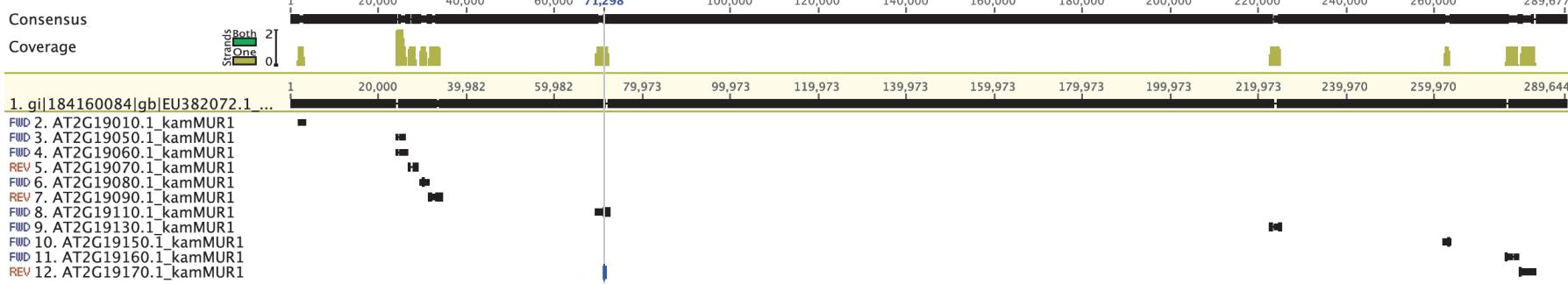


Hanikenne et al. 2013  
PLoS Genetics

**B**



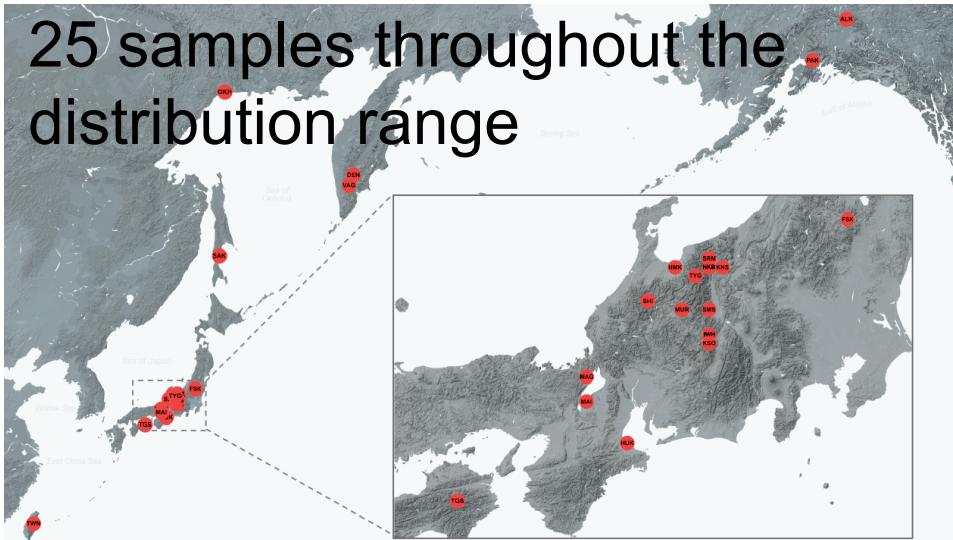
“Hard” selective sweep



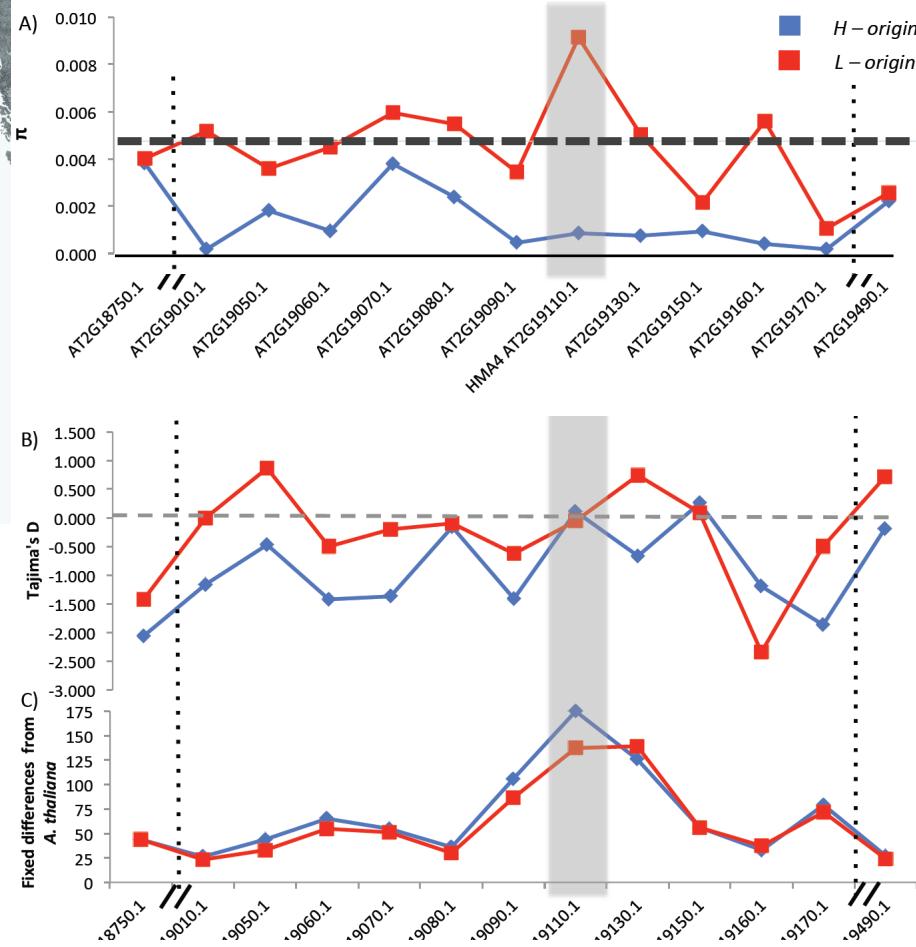
# **Resequencing of 25 genotypes: Signature of selective sweep at the *HMA4* region**

Genome-wide silent nucleotide diversity  $\pi$  = 0.0045 (similar to diploid *A. thaliana*)

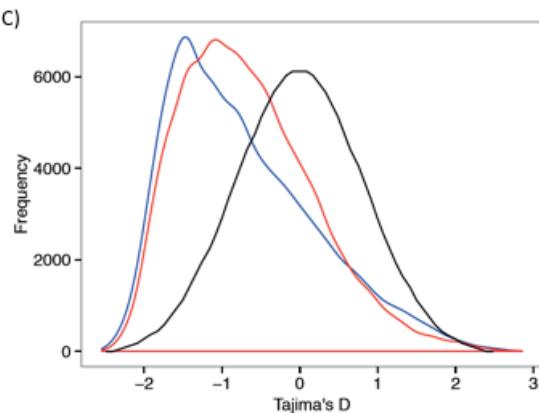
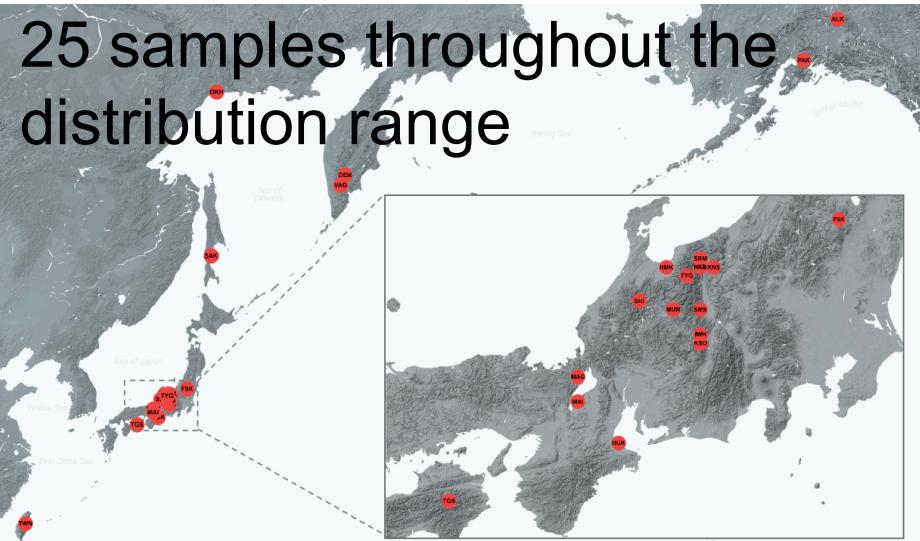
25 samples throughout the distribution range



## Inherited low $\pi$ in *halleri*-derived HMA4 regions



# Different signature of selection



Synonymous nucleotide diversity: correlation is low ( $r = 0.04$ )  
Tajima's  $D$ : Low correlation ( $r = 0.05$ )

Genome-wide pattern:  
different signatures of selection

Increased number of selection targets

# How about other allopolyploid species? Rapid advance in polyploid genome assembly

Hexaploid wheat (17 Gb)  
N50 < 10 kb  
(IWGSC, Science 2015)

>100X  


Finger millet (1.5 Gb)  
N50 ~ 24 kb  
(Hittalmani et al. 2017)

*Brassica juncea*  
N50 ~ 1.5 Mb (Yang et al. 2016)

Tetraploid wheat (12 Gb)  
N50 ~ 7 Mb (Avni et al. 2017)

Hexaploid wheat (Chinese Spring)  
N50 ~ 22.8 MB (2018)

Finger millet (1.5 Gb) Bionano  
N50 ~ 2.6 Mb  
Hatakeyama et al. 2017



La céréale africaine éleusine décodé - SWI swissinfo.ch  
<http://www.swissinfo.ch/fr/g%C3%A9nome-de-la-c%C3%A9r%C3%A9ale-africaine-eleusine-decod%C3%A9.html>

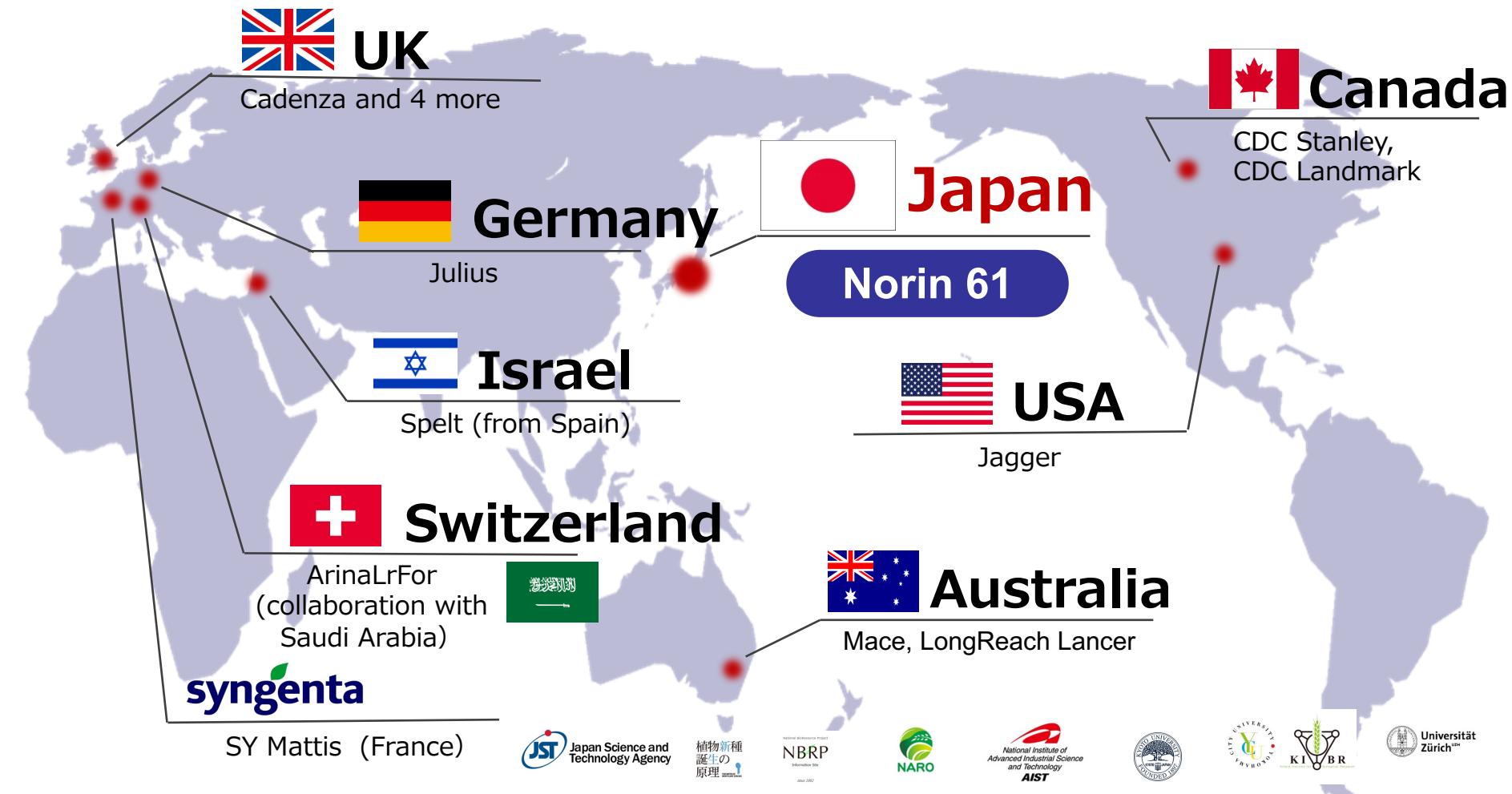
**Genome de la céréale  
africaine éleusine décodé**

 Toute l'actu en bref

05. SEPTEMBRE 2017 - 17:10



# Wheat 10+ genomes project



**Japanese cultivar as a sole Asian representative**

# Wheat 10+ genomes project

Walkowiak et al. *Nature* 288: 577, 2020

Shimizu et al. *PCP* online 2020 (draft cover picture)

## Article

### Multiple wheat genomes reveal global variation in modern breeding

<https://doi.org/10.1038/s41586-020-2961-x>

Received: 3 April 2020

Accepted: 9 September 2020

Published online: 25 November 2020

Open access

Check for updates

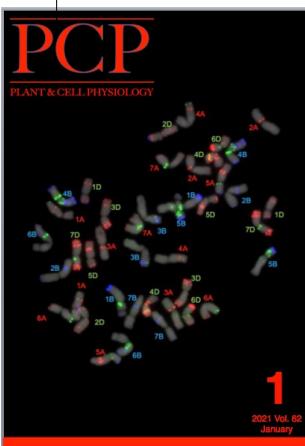
Sean Walkowiak<sup>1,2,4†</sup>, Liangliang Gao<sup>3,4†</sup>, Cecile Monat<sup>4,4†</sup>, Georg Haberer<sup>5</sup>, Mualem T. Kassa<sup>6</sup>, Jemima Brinton<sup>7</sup>, Ricardo H. Ramirez-Gonzalez<sup>7</sup>, Markus C. Kolodziej<sup>8</sup>, Emily Delorean<sup>3</sup>, Dinushka Thambugala<sup>9</sup>, Valentynta Klymiuk<sup>1</sup>, Brook Byrns<sup>1</sup>, Heidrun Gundlach<sup>1</sup>, Venkat Bandi<sup>10</sup>, Jorge Nunez-Sir<sup>10</sup>, Kirby Nilsen<sup>11</sup>, Catharine Aquino<sup>12</sup>, Axel Himmelbach<sup>1</sup>, Dario Copetti<sup>13,14</sup>, Tomohiro Ban<sup>15</sup>, Luca Venturini<sup>16</sup>, Michael Bevan<sup>7</sup>, Bernardo Clavijo<sup>17</sup>, Dal-Hoe Koo<sup>1</sup>, Jennifer Ens<sup>1</sup>, Krystalee Wiebe<sup>1</sup>, Amidou N'Diaye<sup>1</sup>, Allen K. Fritz<sup>2</sup>, Carl Gutwin<sup>18</sup>, Anne Fiebig<sup>19</sup>, Christine Fosker<sup>19</sup>, Bin Xiao Fu<sup>2</sup>, Gonzalo Garcia Accinelli<sup>17</sup>, Keith A. Gardner<sup>20</sup>, Nick Fradley<sup>18</sup>, Juan Gutierrez-Gonzalez<sup>19</sup>, Gwyneth Hallstead-Nussloch<sup>21</sup>, Masaomi Hatakeyama<sup>12,22</sup>, Chi Shih Koh<sup>23</sup>, Jasline Deek<sup>21</sup>, Alejandro C. Costamagna<sup>22</sup>, Pierre Robert<sup>2</sup>, Darren Heavens<sup>7</sup>, Hiroyuki Kanamori<sup>22</sup>, Kanako Kawaura<sup>22</sup>, Fuminori Kobayashi<sup>22</sup>, Ksenia Krasileva<sup>17</sup>, Tony Kuo<sup>4,24,25</sup>, Neil McKenzie<sup>7</sup>, Kazuki Murata<sup>26</sup>, Yusuke Nabeke<sup>26</sup>, Timothy Paape<sup>13</sup>, Sudharshan Padmarasu<sup>4</sup>, Lawrence Percival-Alwynn<sup>19</sup>, Sateesh Kagale<sup>6</sup>, Uwe Scholz<sup>26</sup>, Jun Sese<sup>26,27</sup>, Philimon Julian<sup>28</sup>, Ravi Singh<sup>28</sup>, Rie Shimizu-Itatsugi<sup>19</sup>, David Swarbreck<sup>17</sup>, James Cockram<sup>26</sup>, Hikmet Budak<sup>29</sup>, Toshiaki Tameshige<sup>15</sup>, Tsuyoshi Tanaka<sup>22</sup>, Hirofumi Tsuji<sup>19</sup>, Jonathan Wright<sup>17</sup>, Jianzhong Wu<sup>23</sup>, Burkhard Steuernagel<sup>1</sup>, Ian Smal<sup>19</sup>, Sylvie Cloutier<sup>19</sup>, Gabriel Keeble-Gagnere<sup>22</sup>, Gary Muehlbauer<sup>20</sup>, Josquin Tibbets<sup>2</sup>, Shuhui Nasuda<sup>20</sup>, Joanna Melonek<sup>20</sup>, Pierre J. Hudl<sup>1</sup>, Andrew G. Sharpe<sup>20</sup>, Matthew Clark<sup>16</sup>, Erik Legg<sup>23</sup>, Arvind Bharti<sup>19</sup>, Peter Langridge<sup>34</sup>, Anthony Hall<sup>11</sup>, Cristobal Uauy<sup>1</sup>, Martin Mascher<sup>4,35</sup>, Simon Grattinger<sup>36</sup>, Hirokazu Handa<sup>2,37</sup>, Kentaro K. Shimizu<sup>13,38</sup>, Assaf Distelfeld<sup>38</sup>, Ken Chalmers<sup>34</sup>, Beat Keller<sup>2</sup>, Klaus F. X. Mayer<sup>3,39</sup>, Jesse Poland<sup>2</sup>, Nils Stein<sup>4,40</sup>, Curt A. McCartney<sup>9,23</sup>, Manuel Spannagi<sup>15,22</sup>, Thomas Wicker<sup>8,22</sup> & Curtis J. Pozniak<sup>1,22</sup>

De Novo Genome Assembly of the Japanese Wheat Cultivar Norin 61 Highlights Functional Variation in Flowering Time and *Fusarium* Resistance Genes in East Asian Genotypes

Running head: Genome assembly of the Japanese wheat Norin 61

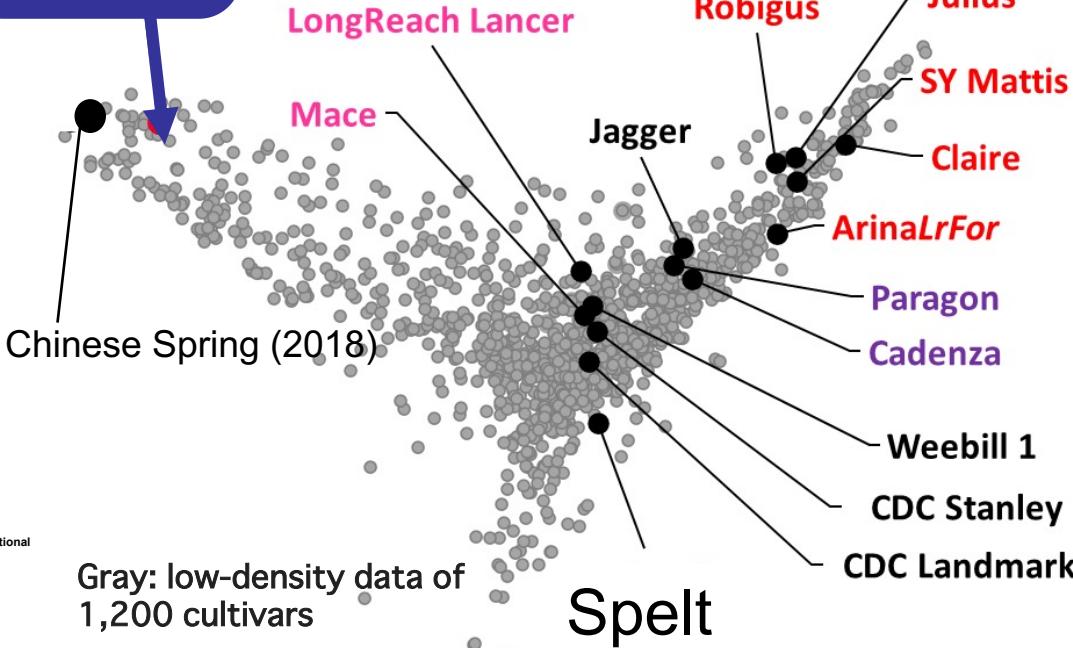
Authors:

Kentaro K. Shimizu<sup>1,2\*#</sup>, Dario Copetti<sup>2,3#</sup>, Moeko Okada<sup>2#</sup>, Thomas Wicker<sup>4</sup>, Toshiaki Tameshige<sup>1,5</sup>, Masaomi Hatakeyama<sup>2,6</sup>, Rie Shimizu-Itatsugi<sup>2</sup>, Catharine Aquino<sup>6</sup>, Kazusa Nishimura<sup>7</sup>, Fuminori Kobayashi<sup>8</sup>, Kazuki Murata<sup>9</sup>, Tony Kuo<sup>10,11</sup>, Emily Delorean<sup>12</sup>, Jesse Poland<sup>12</sup>, Georg Haberer<sup>13</sup>, Manuel Spannagi<sup>11</sup>, Klaus F. X. Mayer<sup>13,14</sup>, Juan Gutierrez-Gonzalez<sup>15</sup>, Gary J. Muehlbauer<sup>15</sup>, Cecile Monat<sup>16</sup>, Axel Himmelbach<sup>16</sup>, Sudharshan Padmarasu<sup>16</sup>, Martin Mascher<sup>16</sup>, Sean Walkowiak<sup>17,18</sup>, Tetsuya Nakazaki<sup>17</sup>, Tomohiro Ban<sup>1</sup>, Kanako Kawaura<sup>1</sup>, Hiroyuki Tsuji<sup>1</sup>, Curtis Pozniak<sup>17</sup>, Nils Stein<sup>16,19</sup>, Jun Sese<sup>1,20\*</sup>, Shuhui Nasuda<sup>20</sup>, Hirokazu Handa<sup>2,21\*</sup>



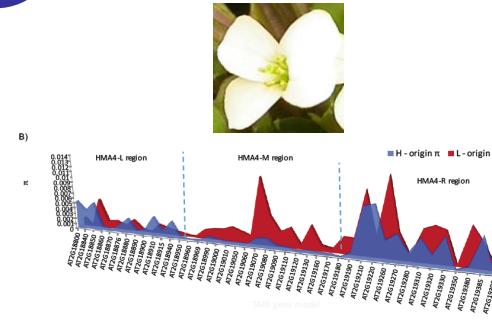
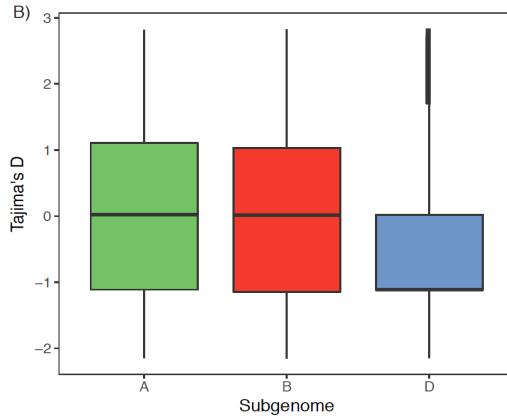
## PCA plot of genome wide polymorphisms

Norin 61



# Different signature of selection among duplicated genes

## Tajima's D (signature of selection)



Paape et al. *Nature Comm*  
9: 3909, 2018

## Low correlation between duplicated copies

Crop  
(bread wheat)

$r = 0.02-0.06$

Natural  
(*Arabidopsis kamchatica*)

$r = 0.05$

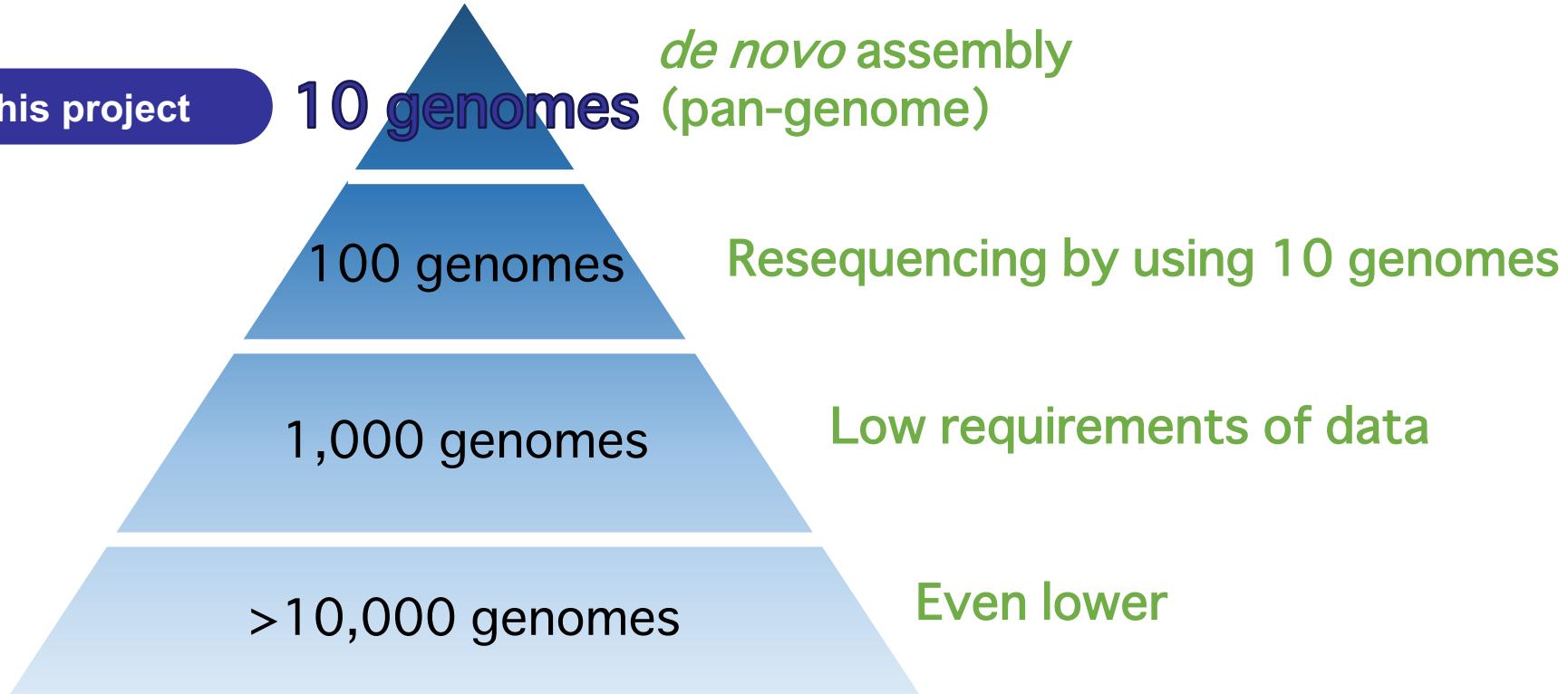
Walkowiak et al. *Nature* 288: 577, 2020

The same pattern in natural and crop polyploid species

# Era of multiple genomes

Experimental strain Chinese Spring  
(*Science*, 2018)

This project



Basic and applied researches