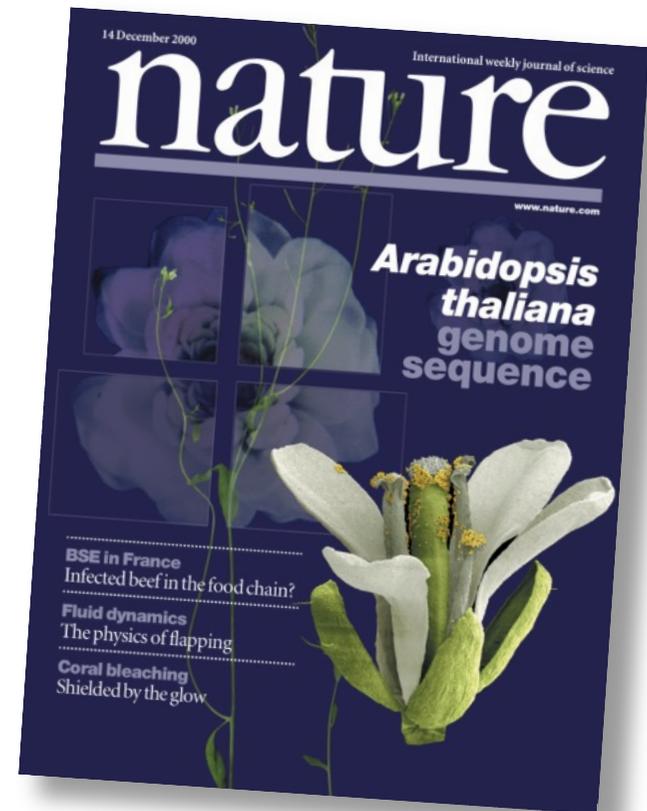# The TAIR12 genome assembly:
## a new reference for *Arabidopsis thaliana*

Xiao Dong, Raúl Wijfjes, Korbinian Schneeberger*
*LMU Munich, Germany*
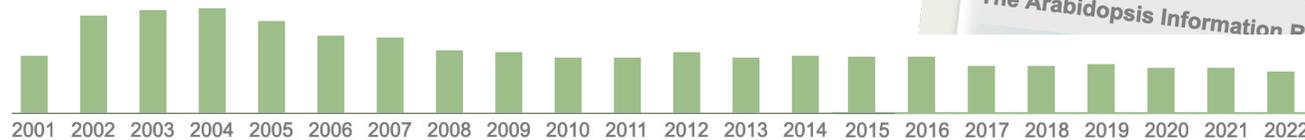*MPI for Plant Breeding Research, Cologne, Germany*

*Community Consensus Arabidopsis thaliana
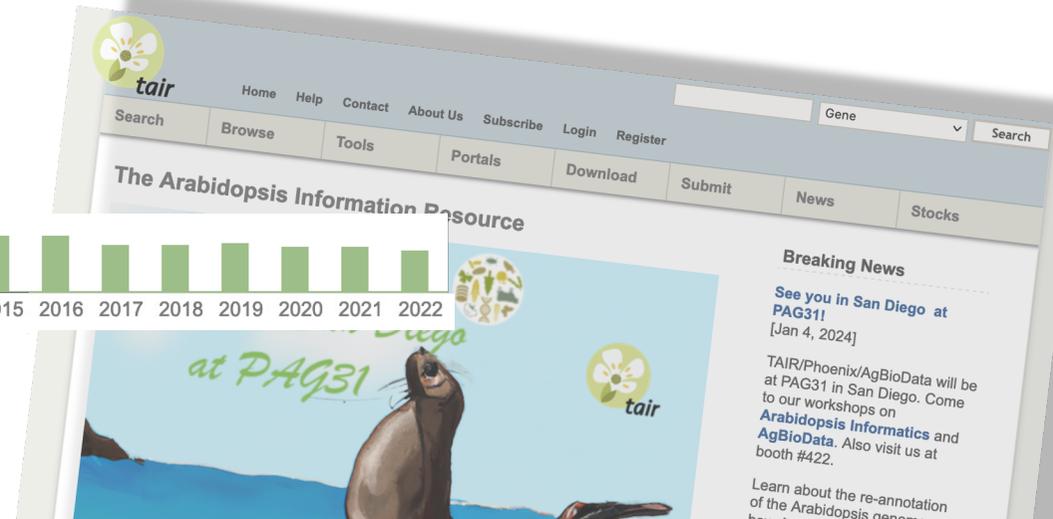Reference Genome Assembly Consortium*

The Arabidopsis Genome Initiative, *Nature*, 2000

# The impact of the *A. thaliana* reference sequences

- Strengthened the model species: Col-0 acts as reference line for functional analysis

- Tool to access the genome: primer and marker design, physical map for mapping

- Acts as community-wide accepted standard for annotations and as knowledge basis (transferred to many other plant genomes)



Google scholar: >11,000 citations

# Long-read sequencing improved genome assembly

Reference sequence not complete (119 Mb vs. ~140 Mb est. genome size)

Naish *et al*, Science, 2021
Wang *et al*, GPB, 2022
Hou *et al*, Mol Plant, 2022
Rabanal *et al*, NAR, 2022

# PacBio HiFi assemblies still include (a few) errors



Assembly 1
Assembly 2
Assembly 3
Assembly 4

Chr1

**Assembly errors**

● Collapse
● Expansion
● Base substitution

· 1-2bp
● 3-50bp
● >50bp

| Contig break

■ Centromere
■ 5S rDNA

4

# PacBio HiFi assemblies still include (a few) errors

# A **community consensus** assembly strategy

Consensus of 13 independent Col-0 assemblies

- (6 different stocks from 5 different labs)



Col-0 community-consensus assembly (**Col-CC**)

Pros / cons of a consensus appraoch

- ❖ No corresponding individual

- ❖ Reconstruction of the ancestral Col-0 genome



Naish et al, 2021; Wang et al, 2022; Hou et al, 2022; Rabanal et al, 2022

# Col-CC.v1: 133 Mb of exceptional quality



7

# The assembly of the nucleolus organizer regions (NORs)



Lysak *et al*, *Methods in Mol Biol*, 2006

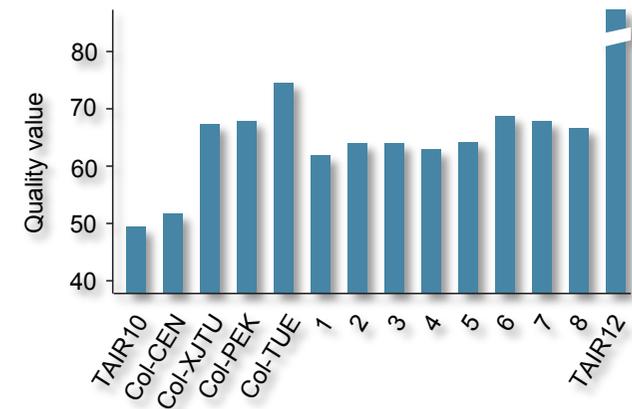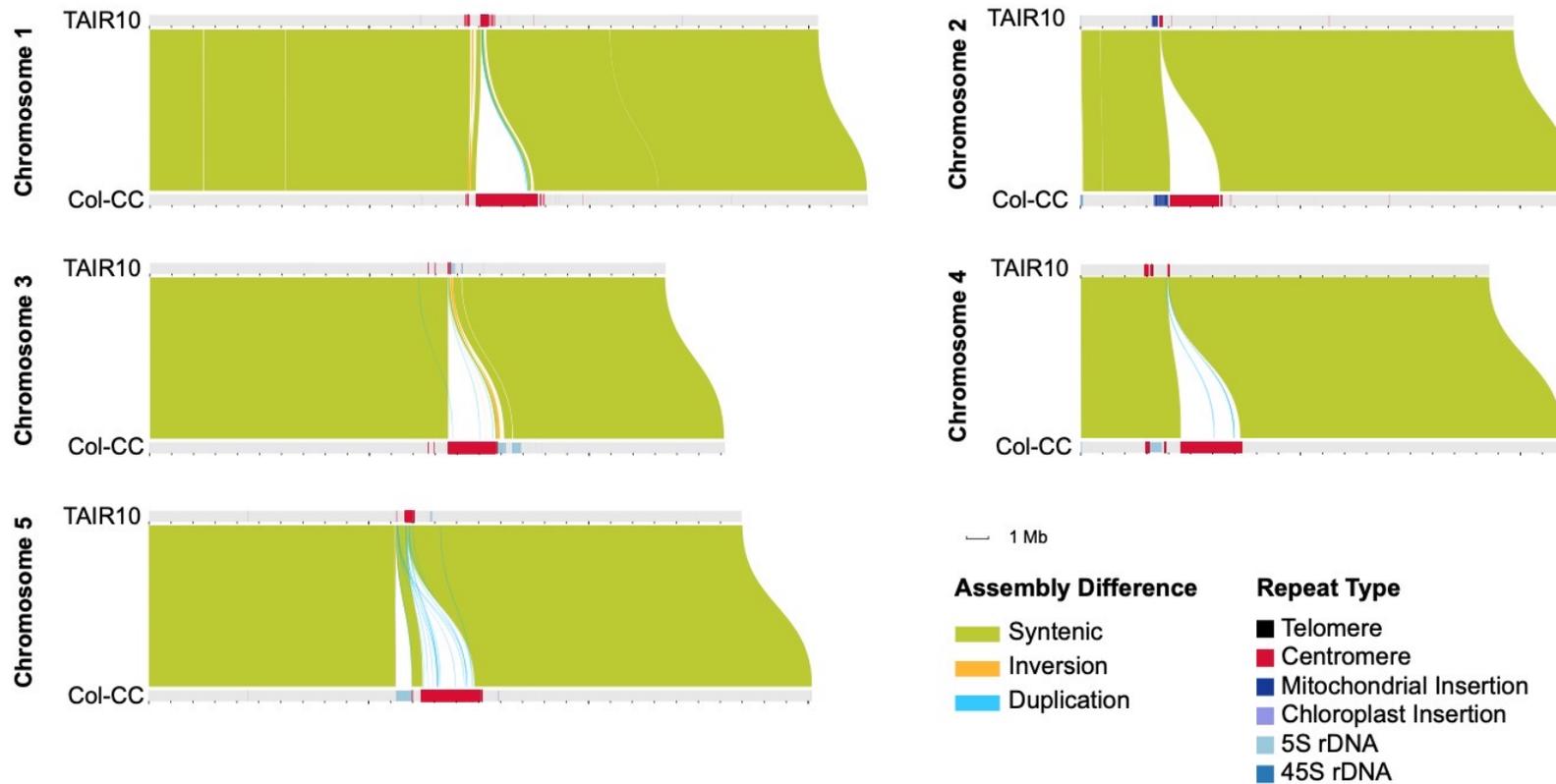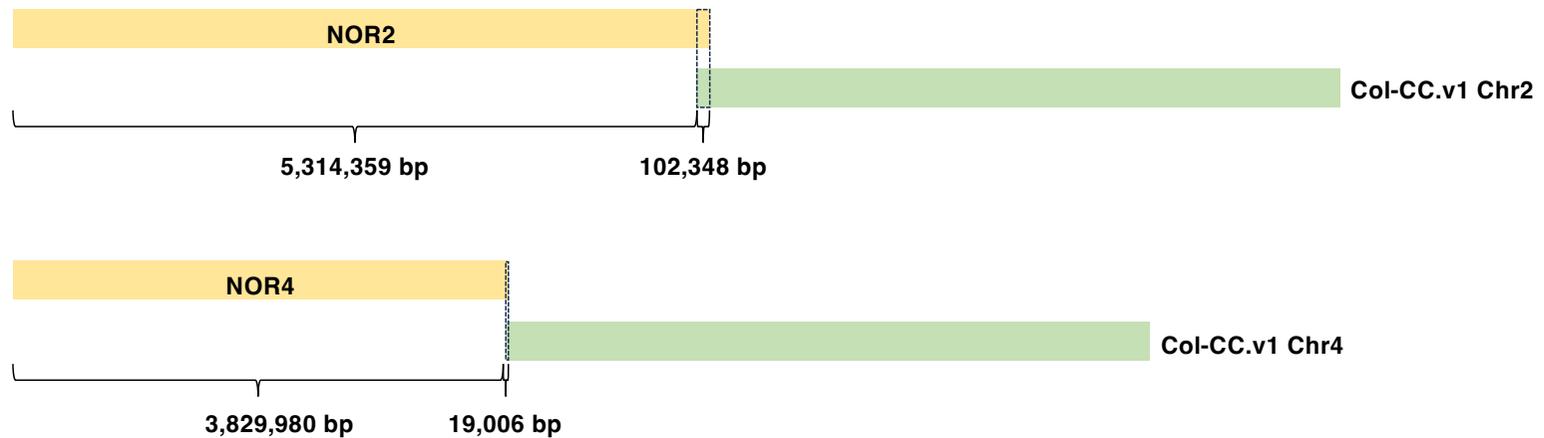Chromosomes 2 and 4: Tel

~3-5 Mb of 45S rDNA repeats

I-*Ppo* I     I-*Ppo* I     I-*Ppo* I     I-*Ppo* I



Computational NOR assembly

ONT reads analyzed for VLE content

Unique VLE landmarks found and connected

Consensus sequence built from all reads containing the unique patterns

Fultz et al, *Science Advances*, 2023

## Assembly of both NOR sequences
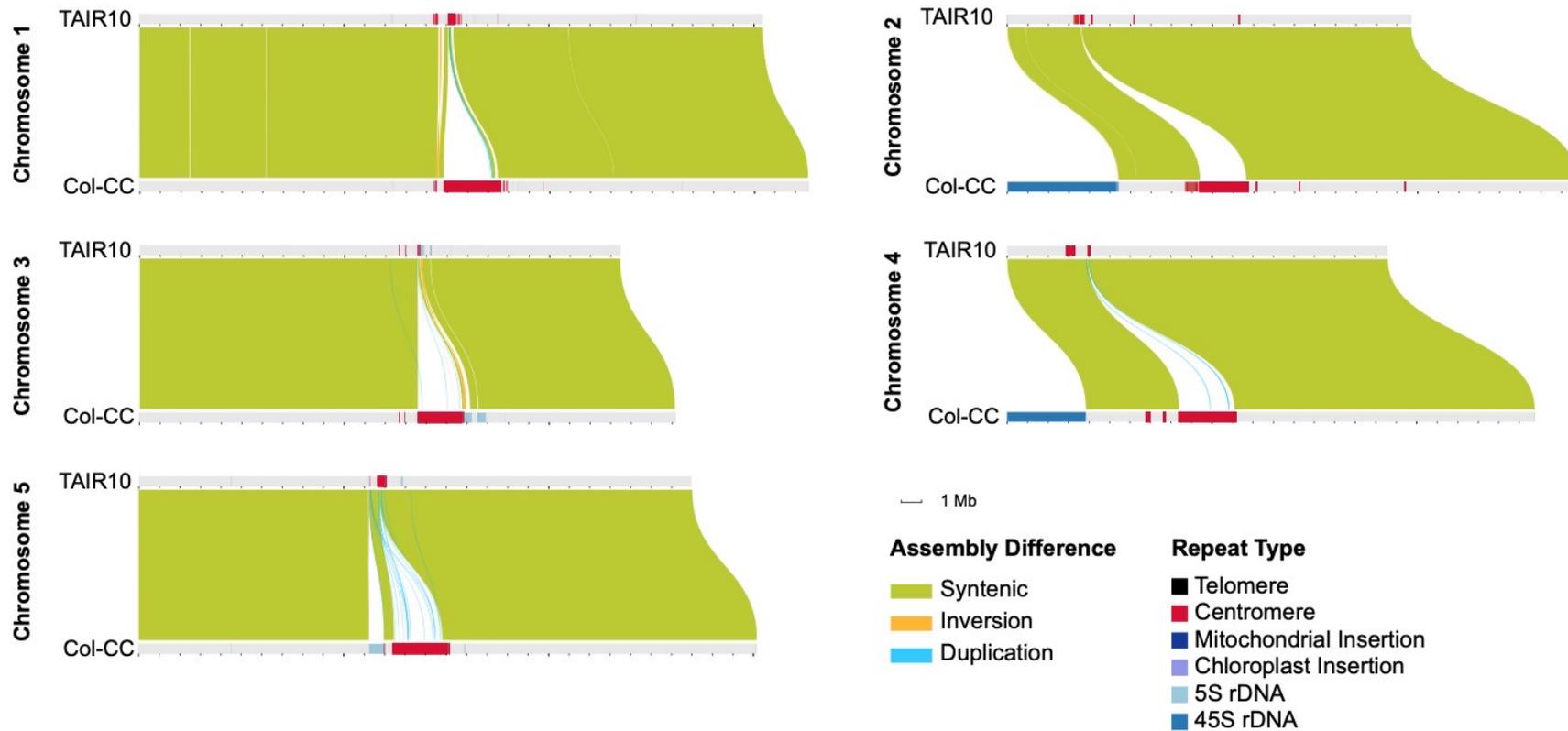
- from telomeres to euchromatin

# Integration of NORs sequences

- NOR2 (5.42 Mb) and NOR4 (3.85 Mb) assemblies overlap with Col-CC

- Extension of Col-CC.v1 at the start of Chr2 and Chr4

# Col-CC.v2: first real gap-less T2T assembly of Arabidopsis

# Col-CC.v2 (assembly of TAIR12): publicly available at NCBI

# Acknowledgement

**Xiao Dong**, Raúl Wijfjes

Henderson, Jiao, Pikaard, Weigel, Ye labs

NCBI    Terence Murphy

tair    **Tanya Berardini**, Leonore Reiser

Annotation groups incl. >70 researchers world-wide contributing to the annotation