*An Introduction into analysis and data generation concepts for complex triticeae genomes – barley&wheat*

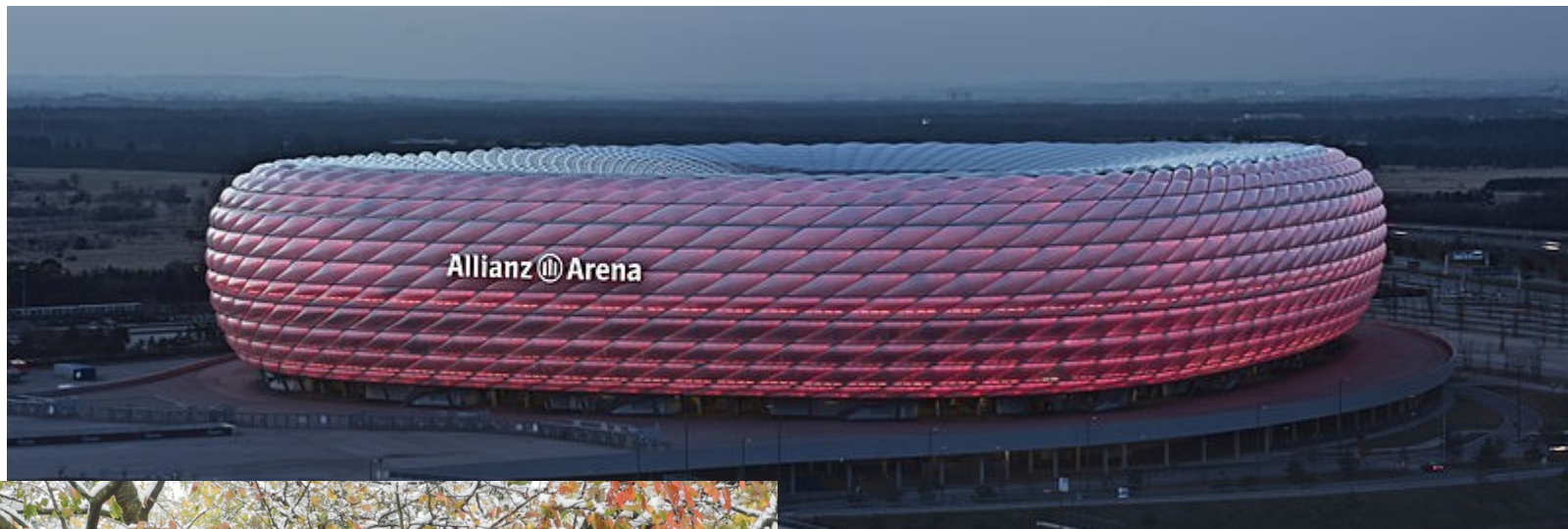*How/where to access the barley&wheat genome data? MIPS PlantsDB tutorial - exercises*

*Manuel Spannagl&Kai Bader*
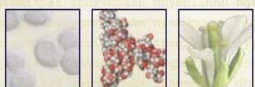*Klaus Mayer*
*MIPS, Helmholtz Center Munich*

transPLANT

HelmholtzZentrum münchen
German Research Center for Environmental Health

CAPACITIES

# Who we are…



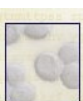c/o Richard Bartz, wikimedia







transPLANT

Helm
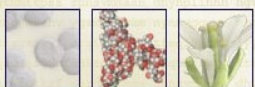German Research Center for Environmental Health

CAPACITIES

# Outline(1)

- Short introduction into data generation and analysis concepts for complex triticeae genomes: **the barley genome**:

  - Barley „genome zipper"
  - Barley genome sequencing, physical+genetic map integration
  - Gene prediction and annotation in barley
  - Comparative genomics in triticeae genomes

# Outline(2)

- Analysing the 17 Gb **genome sequence of bread wheat** using NGS sequencing

- *Optional*: Access to the **Barley physical and genetic maps** – a tutorial intro (based on slides from Nils Stein, IPK)

transPLANT

HelmholtzZentrum münchen
German Research Center for Environmental Health

CAPACITIES

# Outline(3)

- **Barley&Wheat genome database resources**:

  - **MIPS PlantsDB tutorial**: how to access and analyse barley&wheat genome data within a comparative database framework – **interactive exercises**

  - **Homework exercises**: a „real-world" use case accessing the barley genome databases – solutions provided

# BioGreenformatics:
## From Models to Crops, from Pets to Beasts



**Arabidopsis thaliana** *(Nature, 2000)* **Medicago** *(Nature, 2011)* **Tomato** *(Nature, 2012)*

**Sorghum** *(Nature, 2009)* **Brachypodium** *(Nature, 2010)* **Maize** *(Genome Research, 2006, Plant Phys. 2008, PNAS 2008)*

**Barley** *(Plant Phys. 2009, Plant Cell 2011, Nature 2012)* **Arabidopsis lyrata** *(Nature Genetics 2011)*

**Physcomitrella patens** *(Science 2008)* **Aegilops tauschii** *(Nature, under revision)*

**Oryza** *(Genome Research 2001; in prep)* **Wheat** *(Plant Cell 2011;Plant Journal 2012,Nature 2012)*

**Rye** *(PNAS 2012; in prep.)*

**Lolium** *(submitted)* **Festuca** *(submitted)* **Spirodella** *(in prep.)* **Micromonas** *(Science 2009)*

*Cardamine hirsuta*

While NGS democratized sequencing the analytical bottleneck gets more pronounced.

transPLANT

CAPACITIES

# The Challenge

► **Big genome size and high repeat content**



| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 4.6Mb | 157Mb | 300Mb | 430Mb | 730Mb | 3.1Gb | 5.3Gb | 17.1Gb |
| E. coli | Arabidopsis thaliana | Brachy-podium | Rice | Sorghum | Human | Barley | **Bread Wheat** |

# Reduction of complexity by chromosome sorting

## Barley reference sequence - Illumina

Morex 50x WGS assembly 3, repeat masked

 2,670,738 contigs

 1,868,648,155 bp sequence (min 200bp, max 36 kbp, mean 700bp, N50 1,425 bp)

 936,664,164 bp (50.13%) masked sequence

 chromosome arm sorting available (CarmA)

 add. varieties available: Barke, Bowman, …

 SNP/SNVs called and available

# Barley vs *Brachypodium*, *Sorghum* and rice
## -Synteny on a per gene resolution-



[Mb] (window size 0.5 Mb, shift 0.1 Mb)

# Syntenic Integration generates a „GenomeZipper"

# GenomeZipper…what is it? For what?

- Is an approach developed to create an ordered virtual gene map for a chromosome

- It smartly combines chromosome sorting, next generation sequencing, genetic maps, flcDNAs and systematic exploitation of conserved synteny with model grasses

- It provides a valuable surrogate for the gene space of the analyzed chromosome/genome

- Requirements:
  - Masked 454 reads/contigs
  - Orthologs from syntenic regions

# GZipper Input

| marker | flcDNA | brachy | rice | sorghum | reads | ESTs |
|--------|--------|--------|------|---------|-------|------|
| m1 | c1 | b1 | o1 | s1 | r1 | e1 |
| m2 | c2 | b2 | o2 | s2 | r2 | e2 |
| m3 | c3 | b3 | o3 | s3 | r3 | e3 |
| m4 | c4 | b4 | o4 | s4 | r4 | e4 |
| m5 | c5 | b5 | o5 | s5 | r5 | e5 |
| m6 | c6 | b6 | o6 | s6 | r6 | e6 |
| m7 | c7 | b7 | o7 | s7 | r7 | e7 |
| m8 | c8 | b8 | o8 | s8 | r8 | e8 |
| m9 | c9 | b9 | o9 | s9 | r9 | e9 |
| . | . | . | . | . | . | . |
| . | . | . | . | . | . | . |
| mn | cn | bn | on | sn | rn | en |

# Virtual Gene Map: Syntenic Integration



Marker directed synteny projection

Chromosome sorted shotgun sequences

Barley ESTs

High resolution integrated gene map of barley chromosome 1

# GenomeZipper Pipeline

# GenomeZipper pipeline

454 reads/ contigs

→

Modul 1:
Repeat Masking
(Vmatch)

↓

model genomes

marker, flcDNAs, ESTs

→

Modul 2:
Sequence Homology
(BLAST)

↓

Modul 3:
Synteny Detection
(chromoWIZ)

↓

Modul 4:
GenomeZipper

→

Ordered virtual gene map

transPLANT

HelmholtzZentrum münchen
German Research Center for Environmental Health

CAPACITIES

# GenomeZipper: Barley Chromosomes

| Data Sets | 1H mobe | 2H | 3H | 4H | 5H | 6H | 7H | All |
|---|---|---|---|---|---|---|---|---|
| # nonredundant anchored gene loci | 3,331 | 3,616 | 3,394 | 2,709 | 3,208 | 2,304 | 3,204 | 21,766 |
| % markers with associated gene from ref. genome(s) | 63.25 | 61.1 | 66.29 | 69.1 | 60.77 | 58.75 | 53.9 | 61.72 |
| # matched barley fl-cDNAs | 1,676 | 1,619 | 1,628 | 1,255 | 1,474 | 1,058 | 1,395 | 10,105 |
| # nonredundant sequence reads & array hybridization probes | 52,704 | 31,294 | 32,078 | 22,644 | 27,197 | 20,943 | 24,423 | 211,283 |
| # nonredundant ESTs | 3,543 | 3,678 | 3,392 | 2,605 | 3,354 | 2,387 | 3,120 | 22,079 |
| # Brachypodium genes | 2,141 | 2,379 | 2,363 | 1,876 | 2,159 | 1,588 | 1,915 | 14,421 |
| # rice genes | 1,845 | 2,073 | 2,016 | 1,614 | 1,576 | 1,348 | 1,621 | 12,093 |
| # sorghum genes | 1,833 | 1,946 | 2,039 | 1,284 | 1,695 | 1,369 | 1,721 | 11,887 |

transPLANT

HelmholtzZentrum münchen
German Research Center for Environmental Health

CAPACITIES

# Barley – a high resolution genome scaffold

| Chr./ Chr.-arm | expected Lander Waterman of high quality sequences | observed marker detection rate (sensitivity) of high quality sequences | specificity |
|---|---|---|---|
| 1H (MoBe) | 86.46% | 98,19 | 88% |
| 2HS | 64,65% | 82,35 | 97,9 |
| 2HL | 79,20% | 86,24 | 97,1 |
| 3HS | 75,34% | 80,58 | 98 |
| 3HL | 83,14% | 85,95 | 96,5 |
| 4HS | 74,08% | 80,55 | 97,9 |
| 4HL | 78,56% | 83,01 | 93,6 |
| 5HS | 83,63% | 90,29 | 97.9 |
| 5HL | 75,83% | 83,03 | 97,6 |
| 6HS | 82,09% | 86,29 | 97,8 |
| 6HL | 80,60% | 86,38 | 97,8 |
| 7HS | 73,29% | 80,97 | 97 |
| 7HL | 71,35% | 84,89 | 98 |

# Barley Genome Zipper summary

- **22k barley genes** sequence tagged **positionally ordered and** in part associated with flcDNA & EST

- **Additional 6k genes with chromosome arm assignment**

- Resolution of appr. 0,05 (0,1) cM; 20 loci (9,3 fl-cDNAs) per cM

- >3000 (14%) genes are located in low/non-recombining regions

- All but 9 ordered and assigned to short and long arm respectively

## Shortcomings:

- Can't resolve small local rearrangements

- Can't position genes that are out of syntenic context

- Pseudogenes, tandem duplicates, ...

# A powerful shortcut towards an ordered gene map of the huge *Triticeae* genomes



- Rye                              (23k genes zipped)
- *Aegilops tauschii*              (22k genes zipped)
- Wheat – IWGSC                   (70k genes zipped)
- Lolium
- Festuca
- ...

# Genome stratification in barley

Combined genetic and physical map build scaffold

+ Sequence enriched via 80x WholeGenomeShotgun
6,200 BACs,
570,000 BacEndSeqs
500,000 genetic markers
250 Gb RNA Seq

=>3,9Gb (76%) anchored
+ 650 Mb (13%) chr. arm associated

# Different hierarchies of feature connection and different layers to start the navigation

markers

Genetic map

Physical map

Sequence anchored sequence contigs

Genomic features (genes)

# Genome stratification in barley cont.

The Barley **Gene-ome**:

A **physical**, **genetic** and

**functional** sequence

assembly

**More on barley physical and genetic map (integration) in seperate presentation!**

# Gene prediction in barley

## Barley RNA-seq data from SCRI

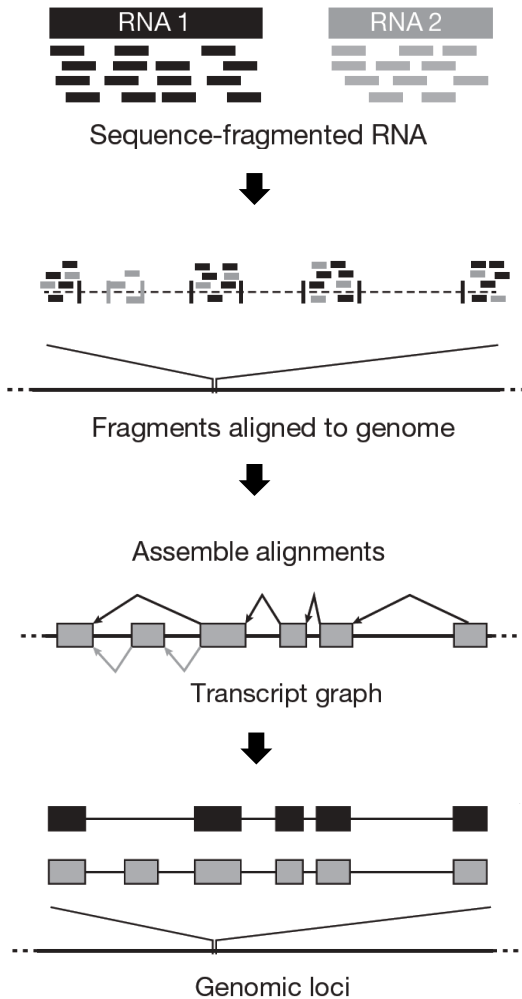| Platform | Read Length [bp] | Paired End | Material | Genotype | Reads [#] | Sequence [bp] |
|---|---|---|---|---|---|---|
| Illumina GA2 SE | 76 | no | 4 days germination embryo | Morex | 23,250,889 | 1,767,067,564 |
| Illumina GA2 SE | 76 | no | 4 days germination embryo | Quench | 26,946,706 | 2,047,949,656 |
| Illumina GA2 SE | 76 | no | 4 days germination embryo | Optic | 23,252,182 | 1,767,165,832 |
| Illumina GA2 SE | 76 | no | 4 days germination embryo | Barke | 25,663,186 | 1,950,402,136 |
| Illumina GA2 SE | 76 | no | 4 days germination embryo | Tocada | 23,868,881 | 1,814,034,956 |
| Illumina GA2 SE | 76 | no | 4 days germination embryo | Betzes | 22,204,022 | 1,687,505,672 |
| Illumina GA2 SE | 76 | no | 4 days germination embryo | Sergeant | 24,480,462 | 1,860,515,112 |
|  |  |  |  |  | Σ 169,666,328 | Σ 12,894,640,928 |

## => Barley reference assembly sequence from IPK

Morex 50x WGS assembly 3, repeat masked

# Gene prediction in barley



Sequence-fragmented RNA

Fragments aligned to genome

Assemble alignments

Transcript graph

Genomic loci

adapted from Garber et al 2011
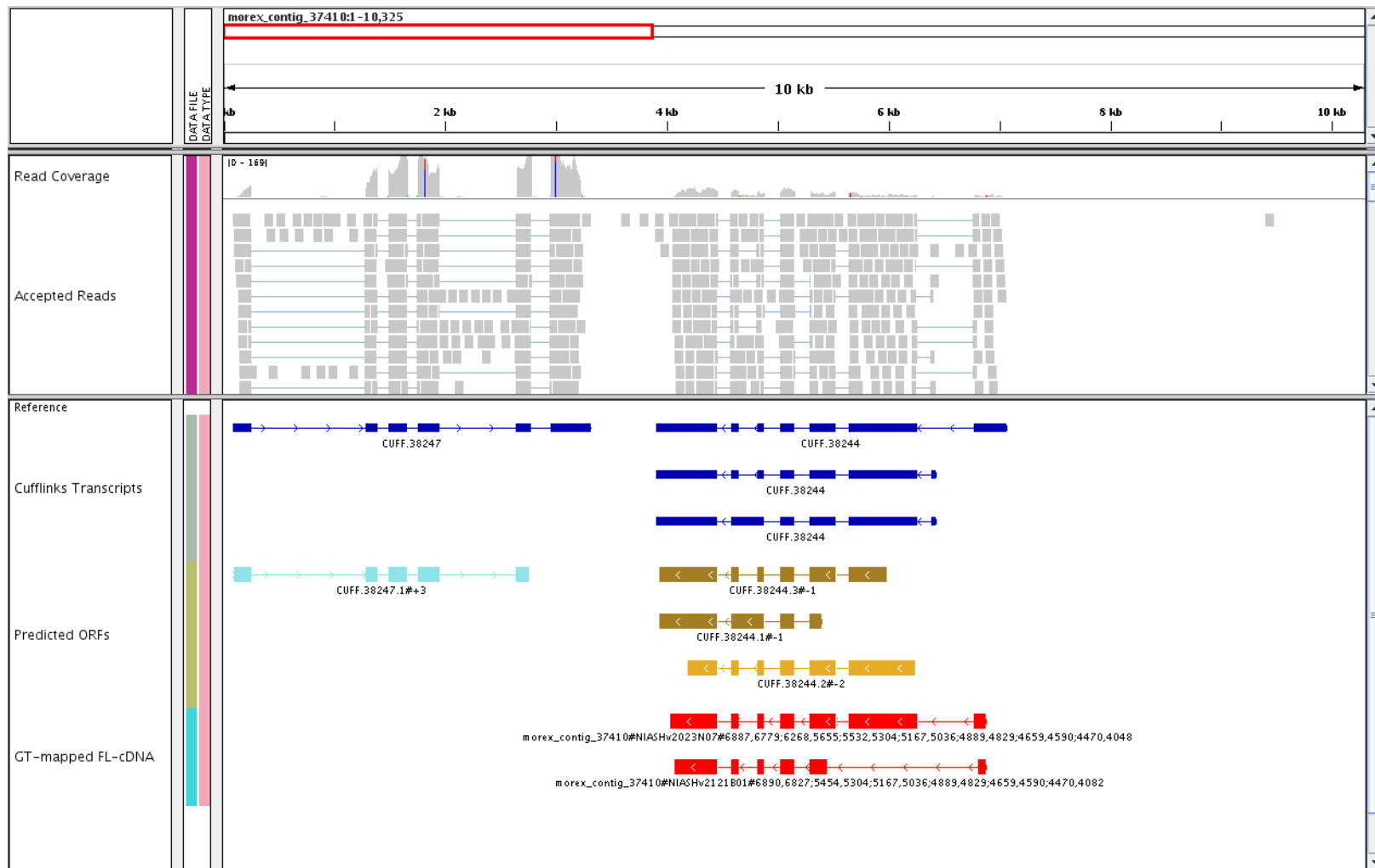
**Bowtie/Tophat**

Alignment of reads to reference genome and identification of splice junctions

**Cufflinks**

Identification of genes and transcripts based on the location of the alignments of spliced reads
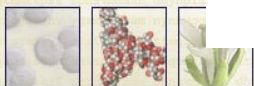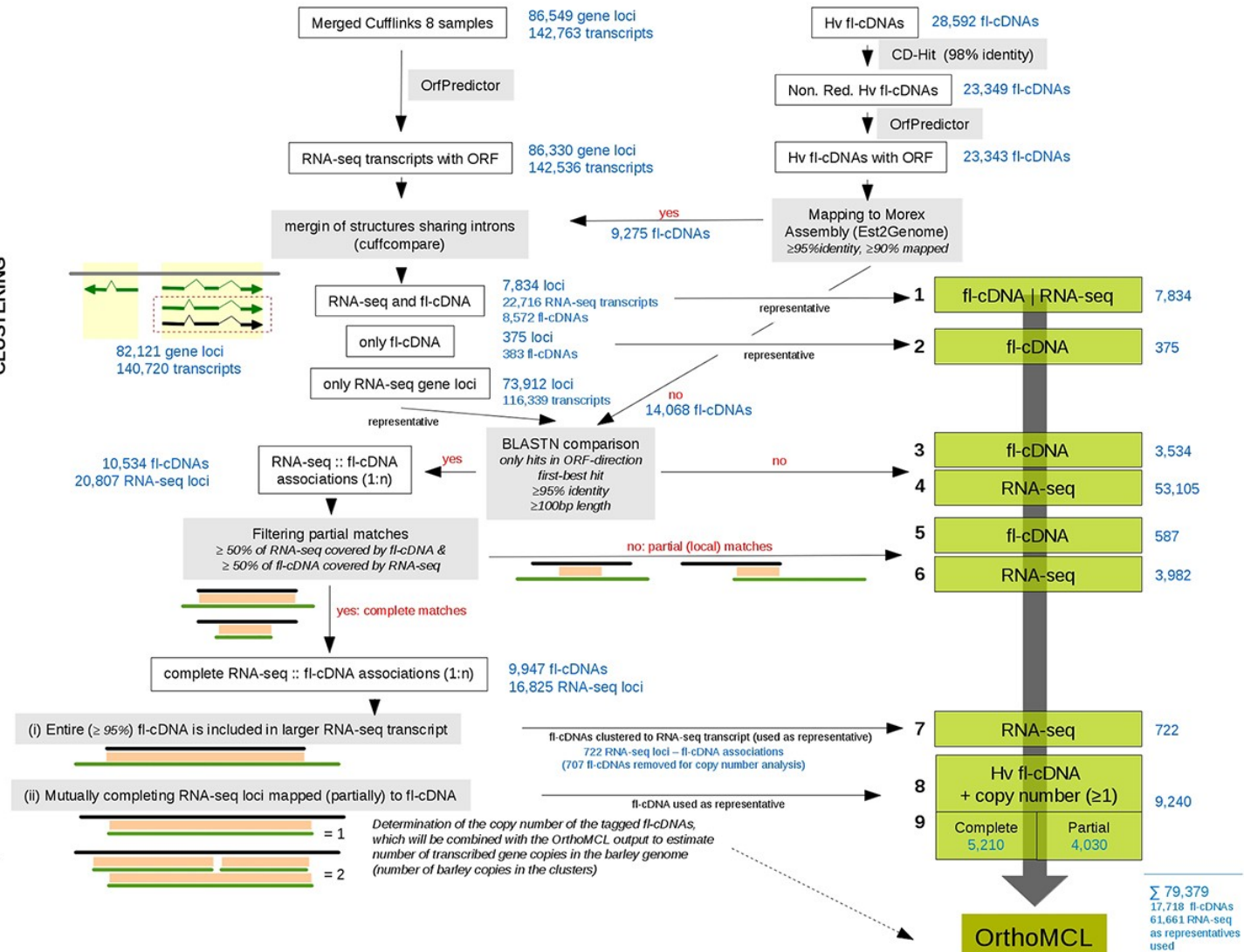
# => 86,330 barley CuffLink loci

# Gene prediction in barley - pipeline
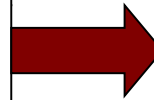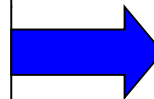
# Gene prediction in barley - results

Total # transcripts clustered for barley (+filtered): **26,159** ➡ "High-confidence" barley genes
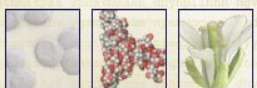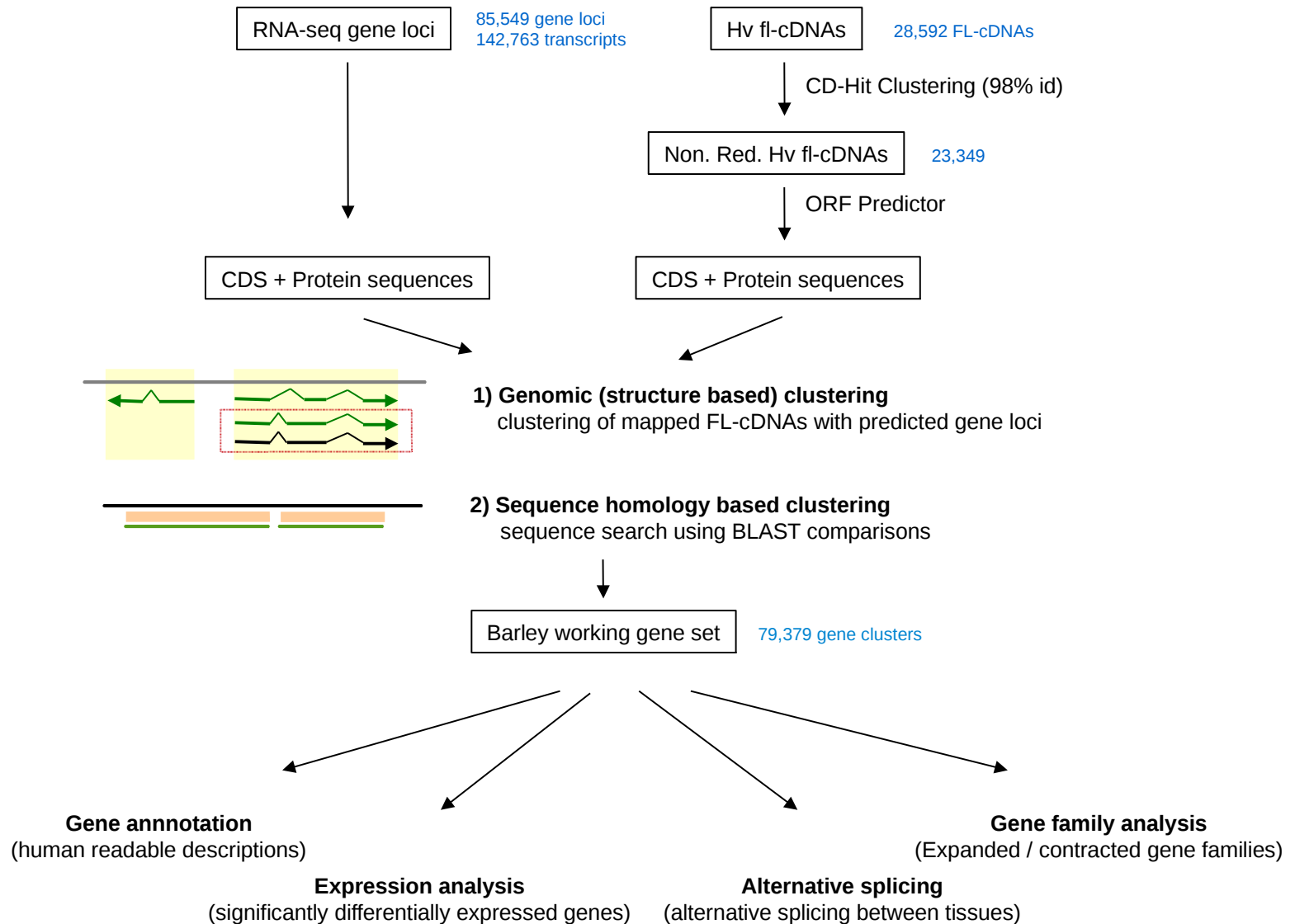
Total # barley Singletons: **53,220** ➡ "Low-confidence" barley genes (likely to contain many pseudogenes & nTARs – novel transcriptional active regions)

# Gene prediction in barley - summary

# acknowledgements