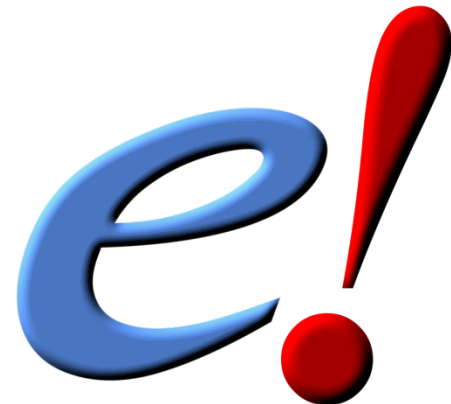# Ensembl Gene Trees & Annotation propagation

Javier Herrero & Glenn Proctor

Vertebrate Genomics Team
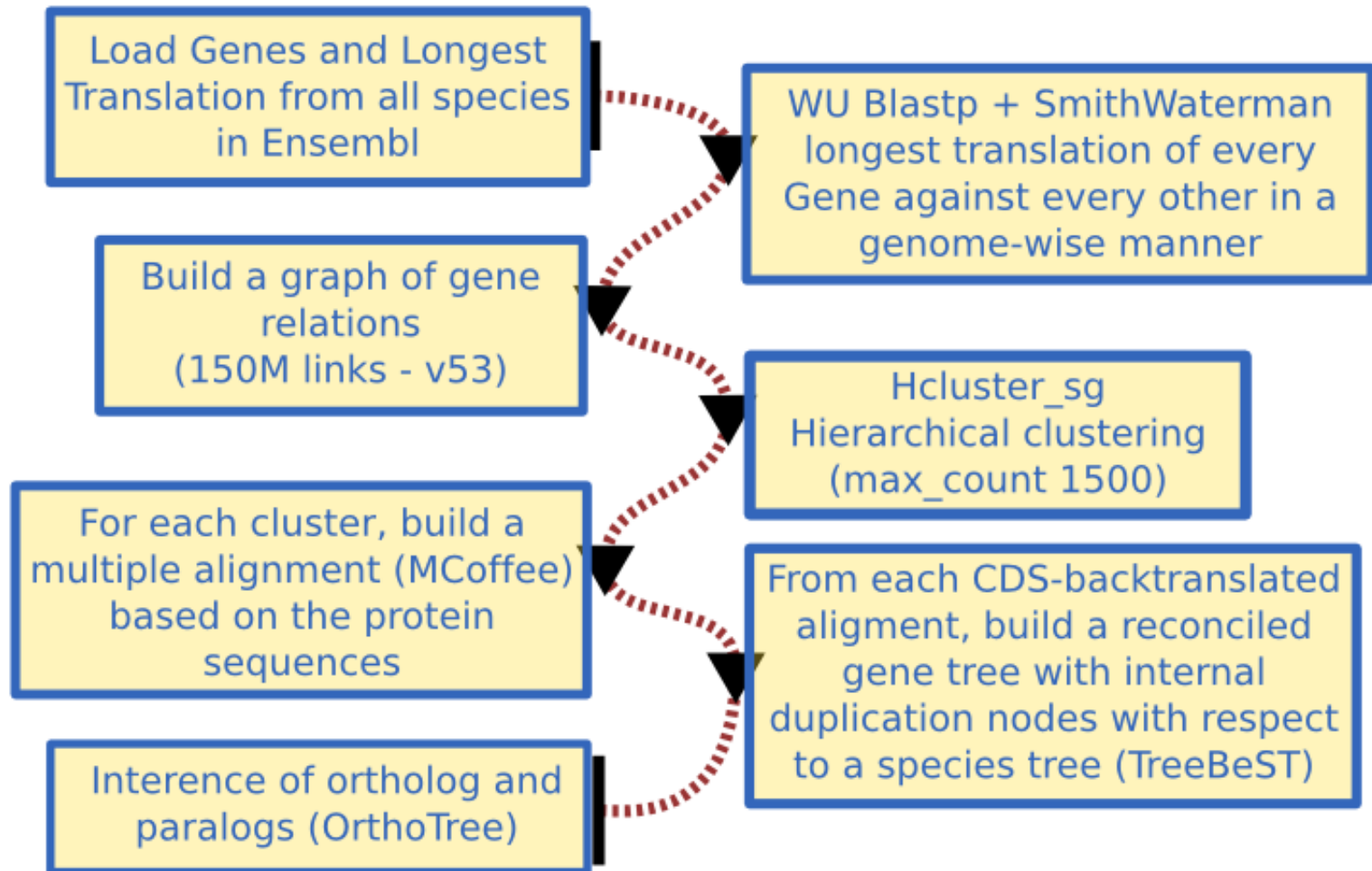
EMBL-EBI

Wellcome Trust Genome Campus

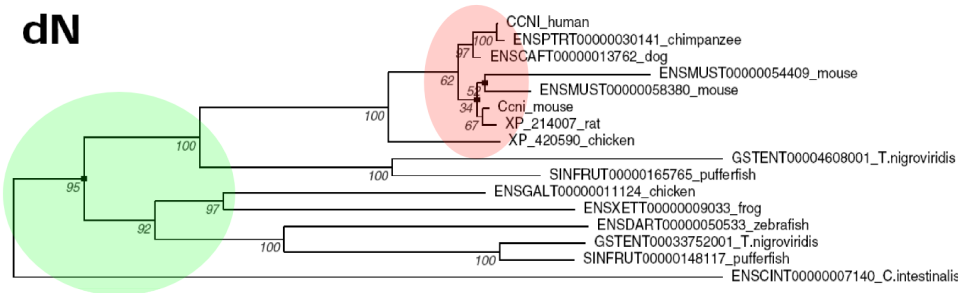Hinxton CB10 1SD, UK

# GeneTree pipeline overview



Vilella et al., Genome Res. 2009

# TreeBeST – treemerge

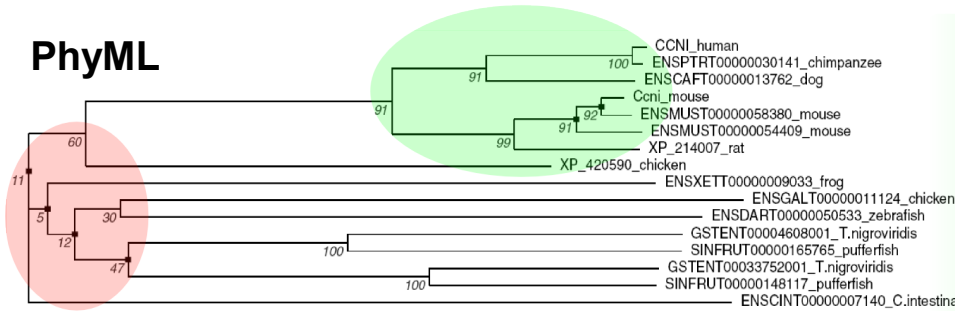- ML-AA-WAG4 – WAG matrix aminoacidic model – Maximum Likelihood (PHYML)

- ML-NT-HKY85 – Hasegawa-Kishino-Yano nucleotidic model – Maximum Likelihood (PHYML)

- NJ-NT-p-distance – any substitutions – neighbor-joining with bootstrap

- NJ-NT-dN – non-syn substitutions – neighbor-joining with bootstrap

- NJ-NT-dS – synonymous substitutions – neighbor-joining with bootstrap

- Curated tree topology (if provided)

# Each method performs best at a given level



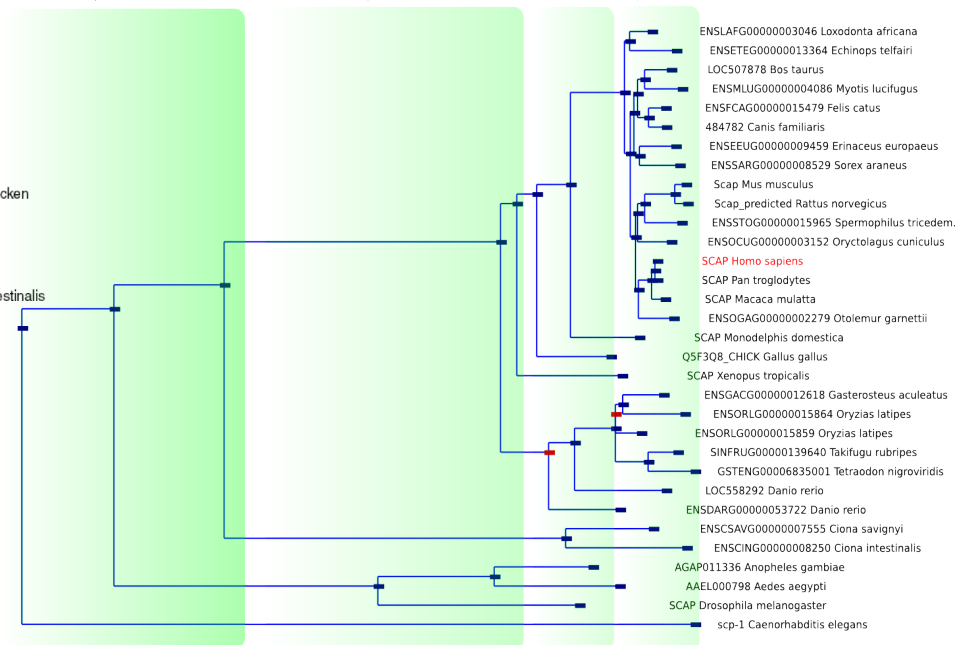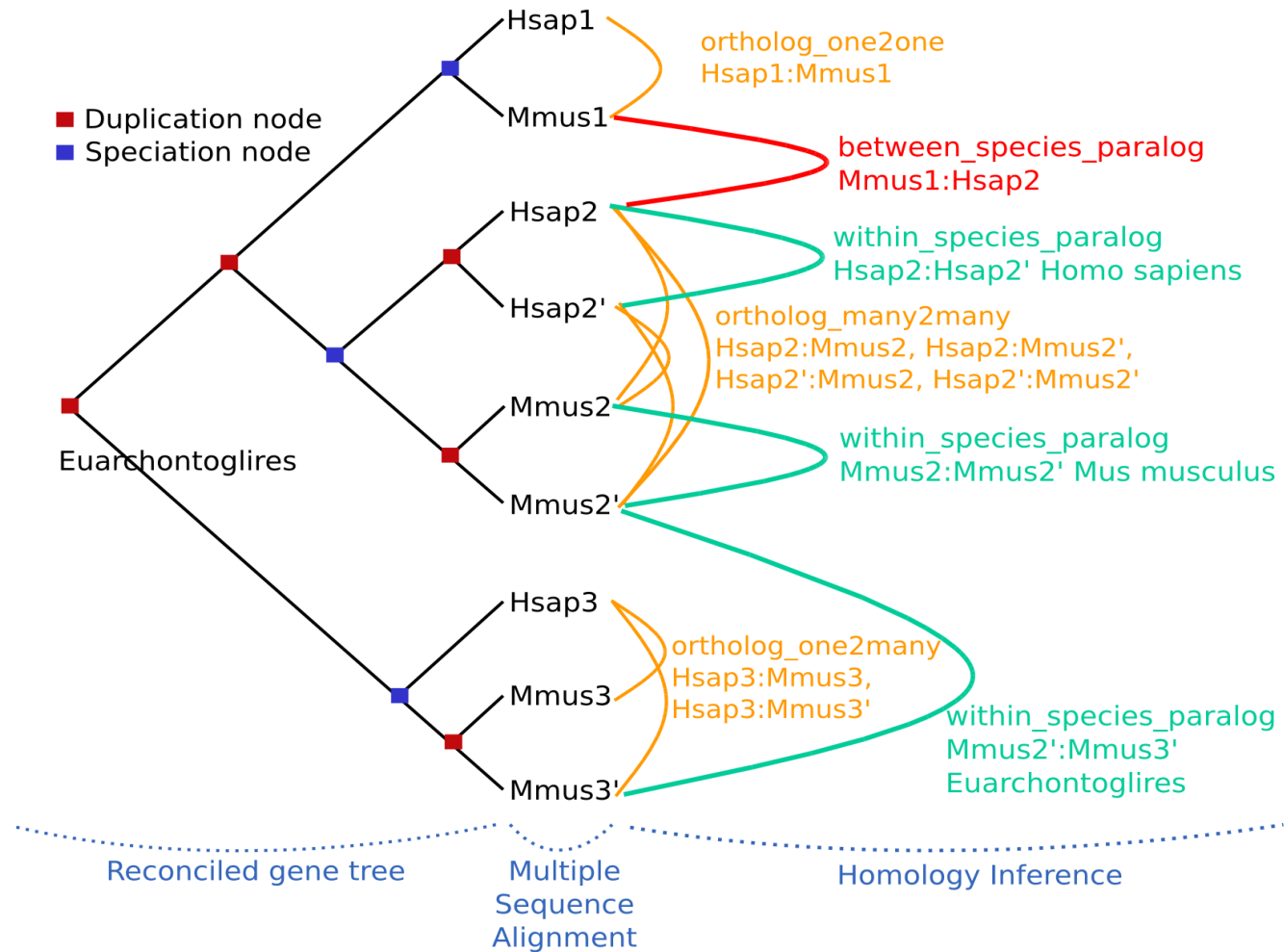**dN**

CCNI_human
ENSPTRT00000030141_chimpanzee
ENSCAFT00000013762_dog
ENSMUST00000054409_mouse
ENSMUST00000058380_mouse
Ccni_mouse
XP_214007_rat
XP_420590_chicken
GSTENT00004608001_T.nigroviridis
SINFRUT00000165765_pufferfish
ENSGALT00000011124_chicken
ENSXETT00000009033_frog
ENSDART00000050533_zebrafish
GSTENT00033752001_T.nigroviridis
SINFRUT00000148117_pufferfish
ENSCINT00000007140_C.intestinalis

**PhyML**

CCNI_human
ENSPTRT00000030141_chimpanzee
ENSCAFT00000013762_dog
Ccni_mouse
ENSMUST00000058380_mouse
ENSMUST00000054409_mouse
XP_214007_rat
XP_420590_chicken
ENSXETT00000009033_frog
ENSGALT00000011124_chicken
ENSDART00000050533_zebrafish
GSTENT00004608001_T.nigroviridis
SINFRUT00000165765_pufferfish
GSTENT00033752001_T.nigroviridis
SINFRUT00000148117_pufferfish
ENSCINT00000007140_C.intestinalis

Method D   Method C   Method B   Method A

ENSLAFG00000003046 Loxodonta africana
ENSETEG00000013364 Echinops telfairi
LOC507878 Bos taurus
ENSMLUG00000004086 Myotis lucifugus
ENSFCAG00000015479 Felis catus
484782 Canis familiaris
ENSEEUG00000009459 Erinaceus europeaus
ENSSARG00000008529 Sorex araneus
Scap Mus musculus
Scap_predicted Rattus norvegicus
ENSSTOG00000015965 Spermophilus tricedem.
ENSOCUG00000003152 Oryctolagus cuniculus
SCAP Homo sapiens
SCAP Pan troglodytes
SCAP Macaca mulatta
ENSOGAG00000002279 Otolemur garnettii
SCAP Monodelphis domestica
Q5F3Q8_CHICK Gallus gallus
SCAP Xenopus tropicalis
ENSGACG00000012618 Gasterosteus aculeatus
ENSORLG00000015864 Oryzias latipes
ENSORLG00000015859 Oryzias latipes
SINFRUG00000139640 Takifugu rubripes
GSTENG00006835001 Tetraodon nigroviridis
LOC558292 Danio rerio
ENSDARG00000053722 Danio rerio
ENSCSAVG00000007555 Ciona savignyi
ENSCING00000008250 Ciona intestinalis
AGAP011336 Anopheles gambiae
AAEL000798 Aedes aegypti
SCAP Drosophila melanogaster
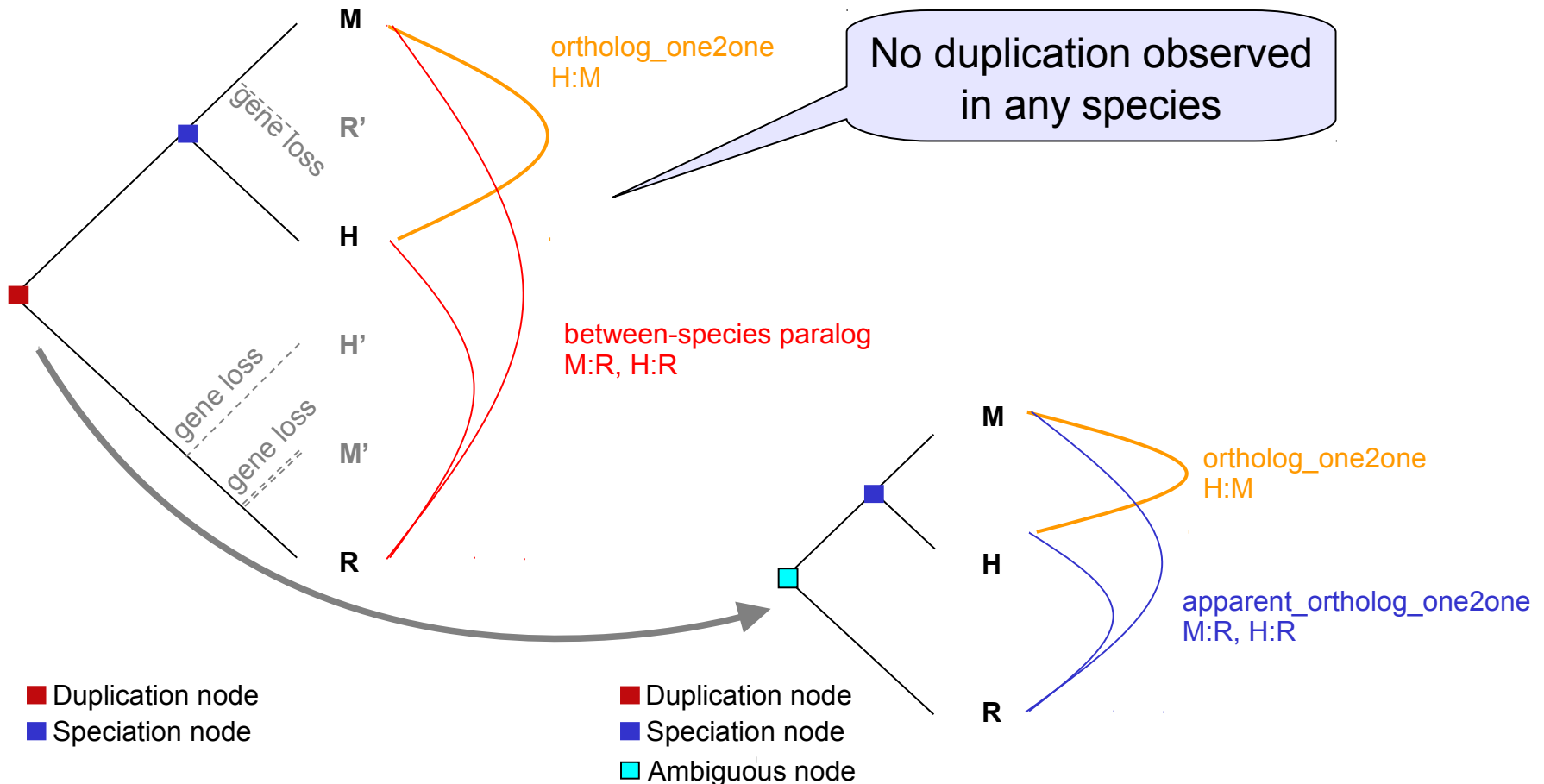scp-1 Caenorhabditis elegans

wellcome trust
sanger
institute

EMBL-EBI
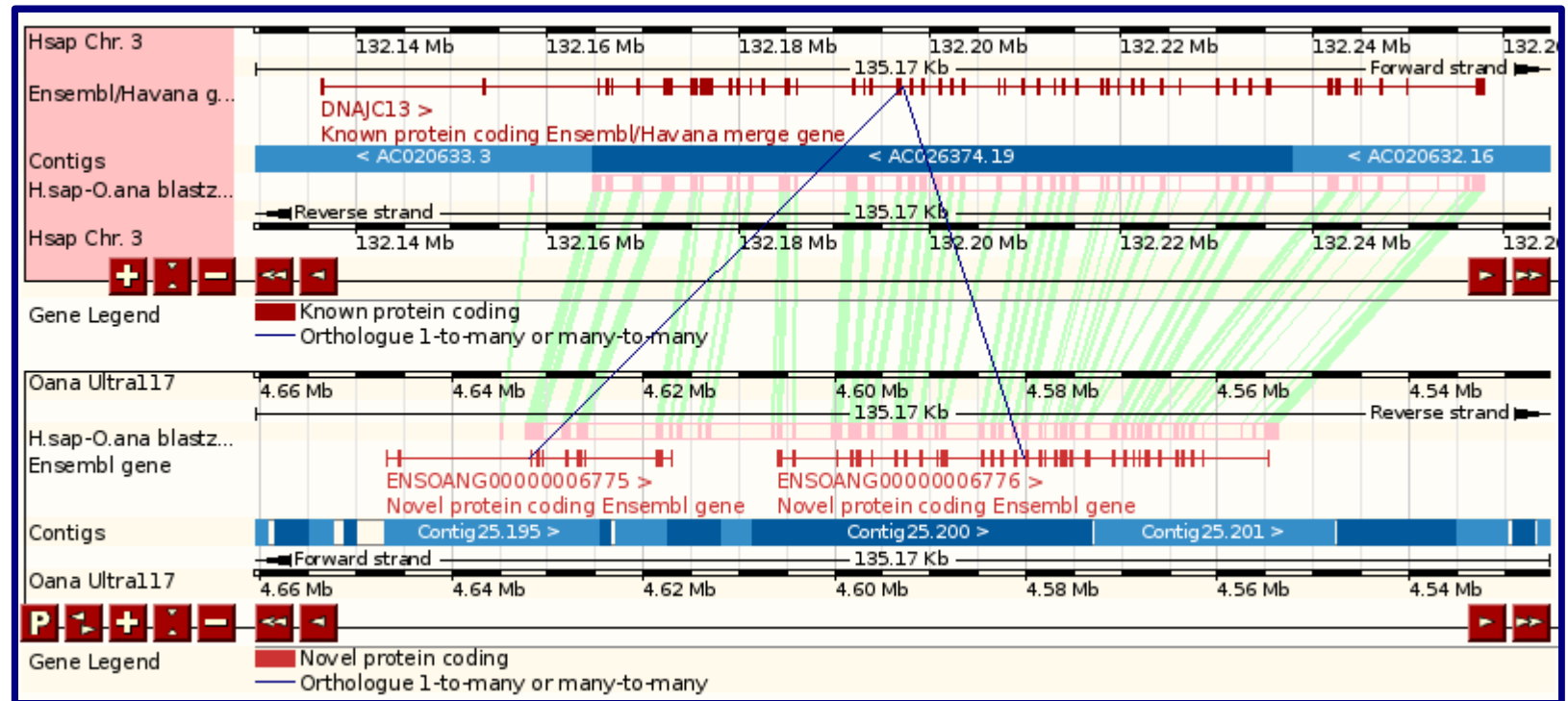
# Homology inference

# A special case of homology



**Orthologues** : any gene pairwise relation where the ancestor node is a SPECIATION event.
**Paralogues** : any gene pairwise relation where the ancestor node is a DUPLICATION event.
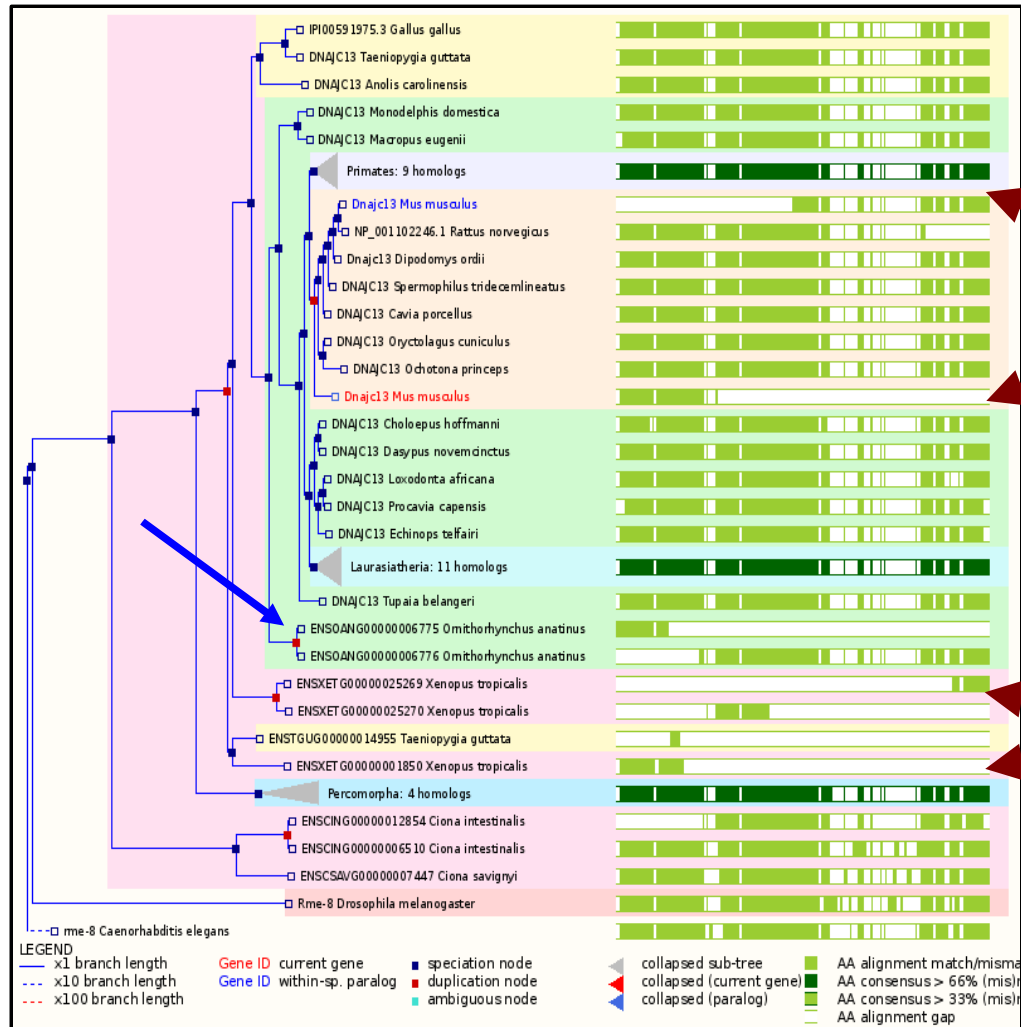
# Gene splits

- DNAJC13 is split in platypus:



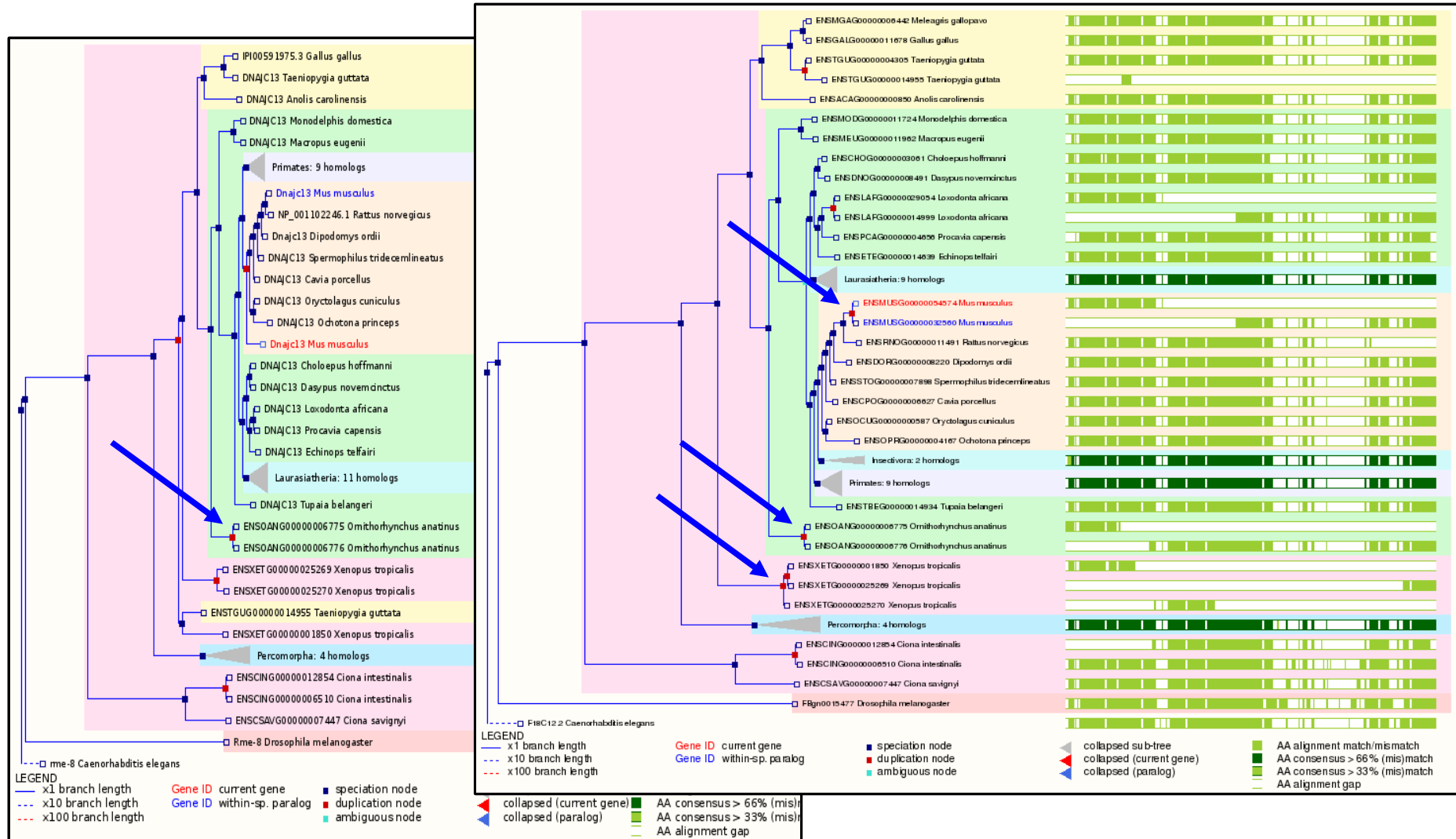- We annotate them as gene split instead of within-species paralogues

# Gene splits



- Most of the splits are in:

  – anole lizard, platypus, zebrafish, macaque, frog, dog, zebra finch

- Feedback to genebuilders

- Ideally, we want to see them together in the tree

- Copy & paste the missing sequence to the other part of the gene

# Gene splits

# Orthologues in forthcoming e! 58

- ### Human-dog

| Orthology type | Pairwise events |
|---|---|
| 1:1 orthologues | 14932 |
| Apparent 1:1 orthologues | 445 |
| 1:many orthologues | 2193 (1520 human genes) |
| many:many orthologues | 2095 (358 human genes) |

- ### Human-zebrafish

| Orthology type | Pairwise events |
|---|---|
| 1:1 orthologues | 8779 |
| Apparent 1:1 orthologues | 187 |
| 1:many orthologues | 8633 (4582 human genes) |
| many:many orthologues | 5034 (933 human genes) |

NB: These values include ncRNA genes as well

# Gene name projections

- From human and mouse to all vertebrates

    - Use the 1:1 orthologs

- From human to fish

    - Use the 1:many orthologs

    - Names become NAME (1 of 3), NAME (2 of 3), etc.

- Rules:

    - if source gene has an HGNC name, and target gene has no name or only a RefSeq predicted name, add name to target gene.

    - Change status to "KNOWN_BY_PROJECTION"

wellcome trust
**sanger**
institute

EMBL-EBI

# GO term projections

- Use the 1:1 orthologs only

- From human and mouse to all vertebrates

- From rat to human and mouse

- From zebrafish to other fish

- From human to zebrafish

- Rules:

  - add GO terms from source gene to target gene, avoiding duplicates. Only project source GO terms with evidence codes IDA, IEP, IGI, IMP, IPI.

  - Projected GO terms on target gene are given evidence code IEA.