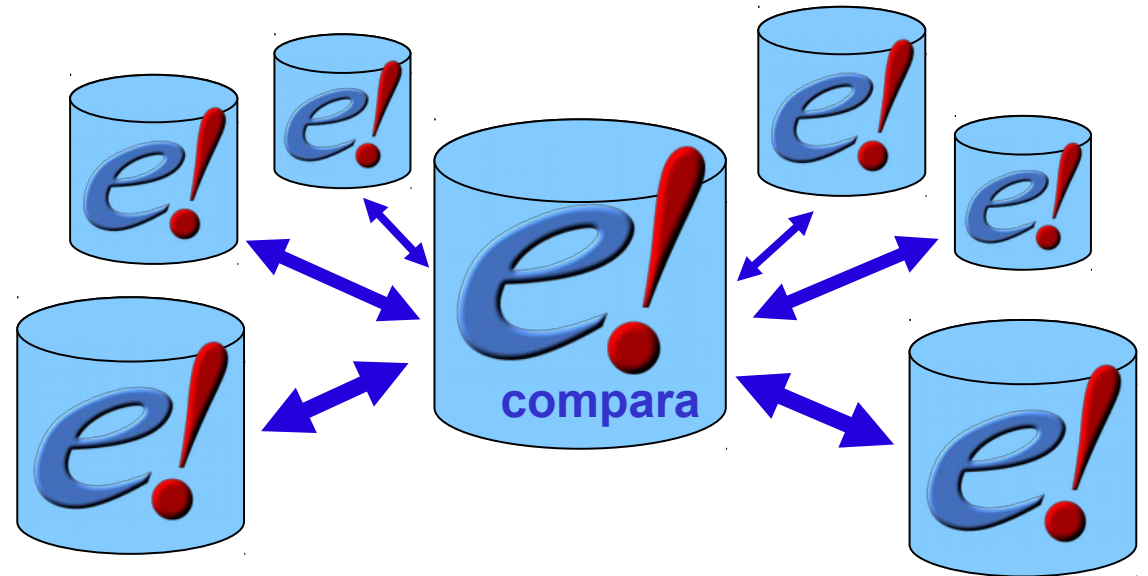
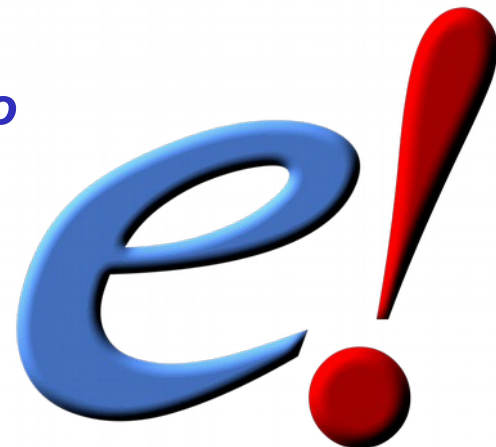


# Ensembl Compara Perl API



*Carla Cummins & Matthieu Muffato*

*API workshop - EBI*



# Outline of the course

- Introduction about Compara
  - Resources
  - API
- Inputs
  - Species, Chromosomes, Genes
- Outputs
  - Gene analyses
  - Genome analyses

# Outline of the course

- Introduction about Compara
  - Resources
  - API
- Inputs
  - Species, Chromosomes, Genes
- Outputs
  - Gene analyses
  - Genome analyses



# What is Ensembl Compara?

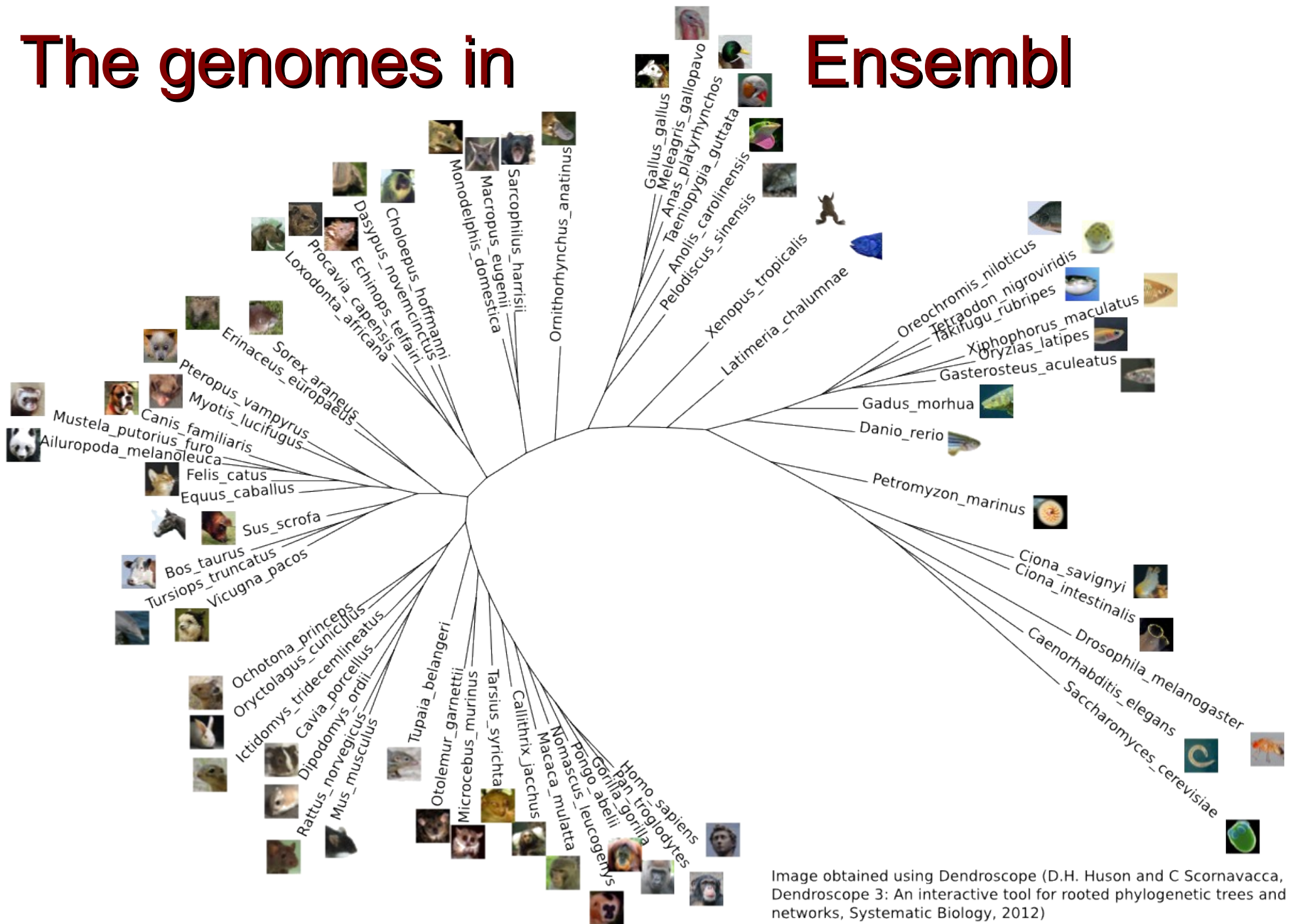
A single database which contains precalculated comparative genomics data and which is linked to all the Ensembl Species (69 in e81) databases.

Access via perl API and mysql

A production system for generating that database  
*(not in this presentation)*

# The genomes in

# Ensembl



# Compara data

## Genome level

Whole genome alignments (pairwise and multiple)

Constrained elements (based on multiple align.)

Syntenic regions (based on pair-wise align.)

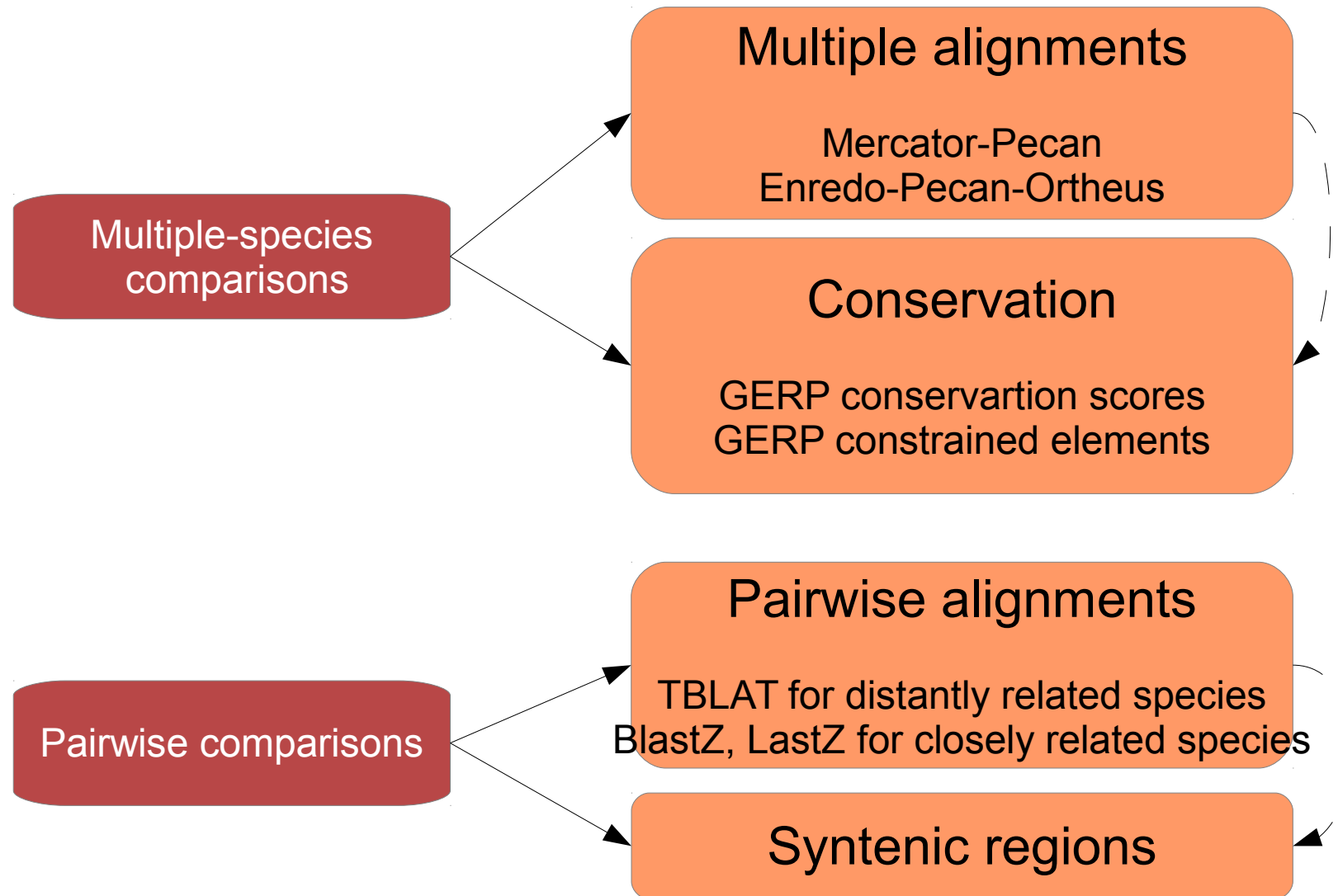
## Gene level

Families (clusters of proteins + multiple align.)

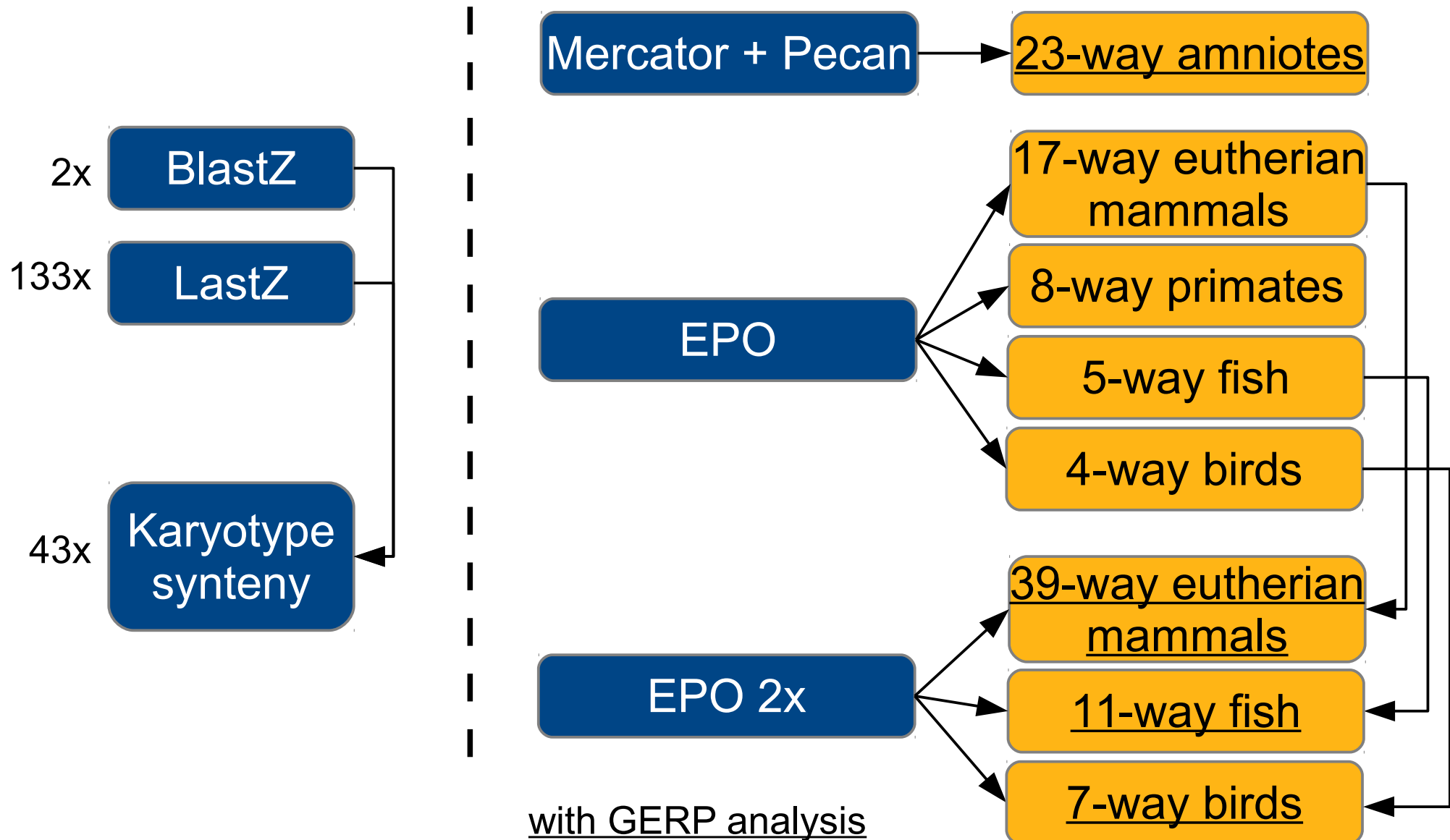
Gene trees (proteins, non-coding RNAs)

Gene orthology / paralogy predictions

# Nucleotide sequence analyses

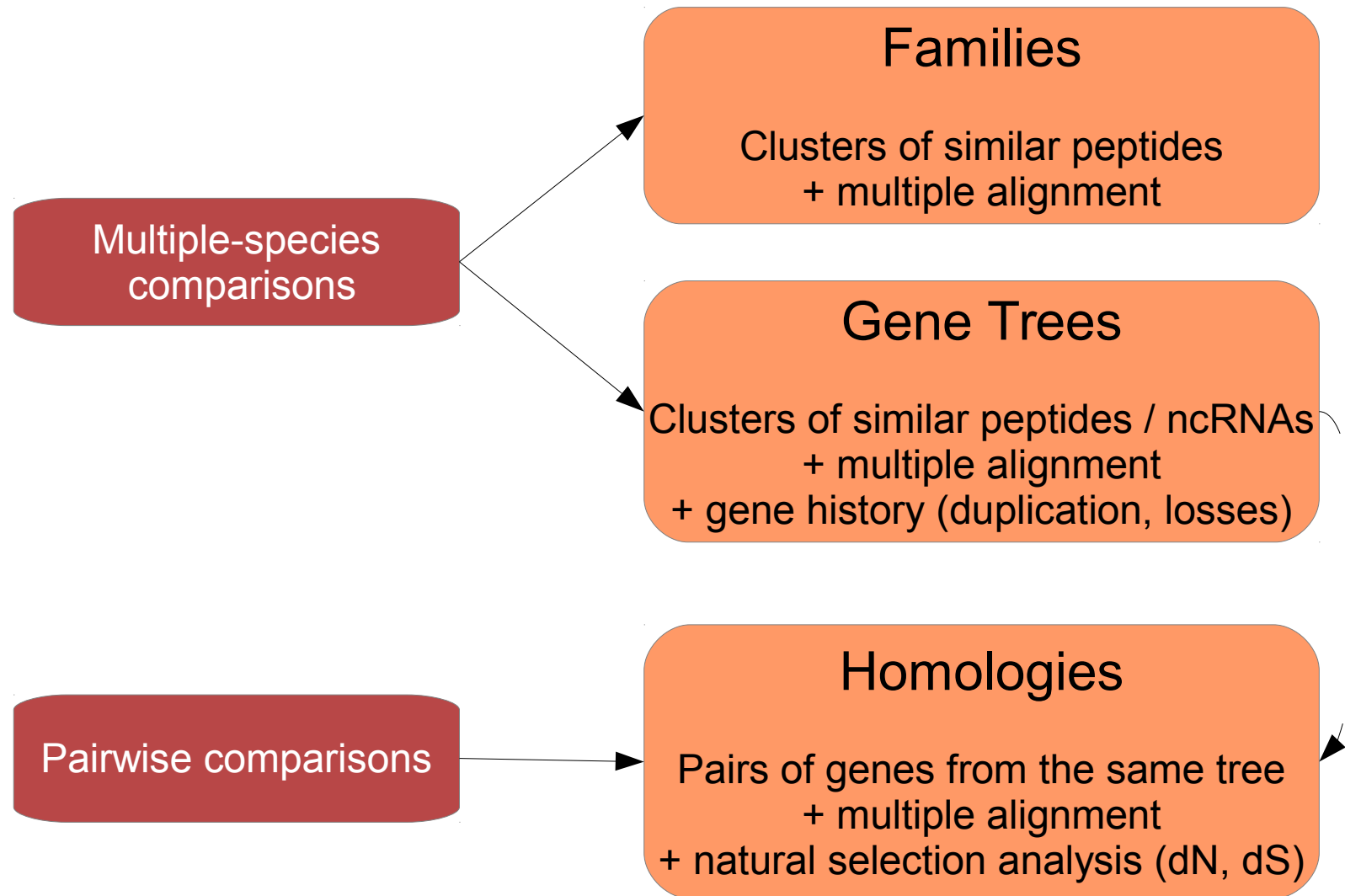


# Nucleotide sequence analyses in e!81





# Gene analyses



# Outline of the course

- Introduction about Compara
  - Resources
  - API
- Inputs
  - Species, Chromosomes, Genes
- Outputs
  - Gene analyses
  - Genome analyses

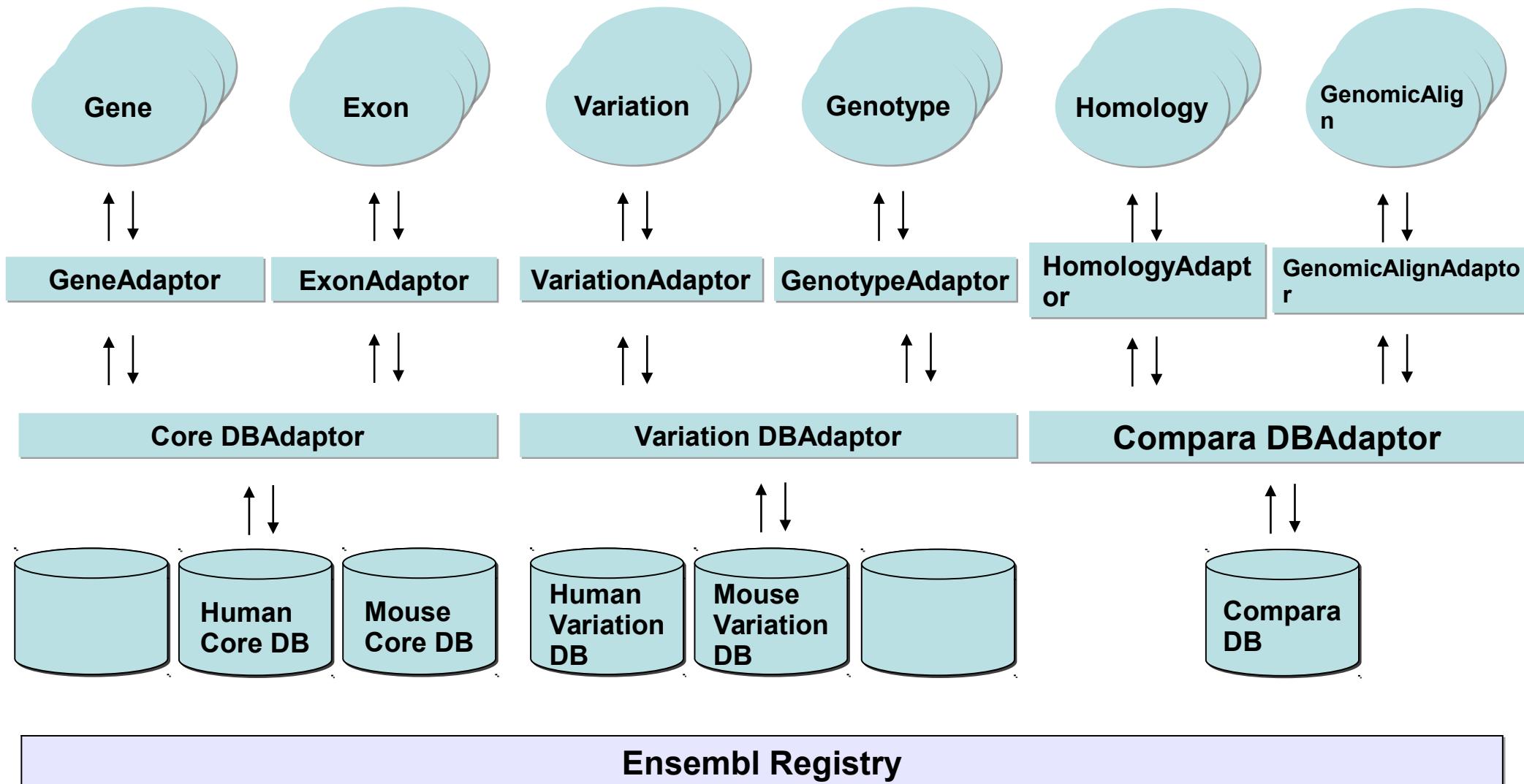


# The Compara Perl API

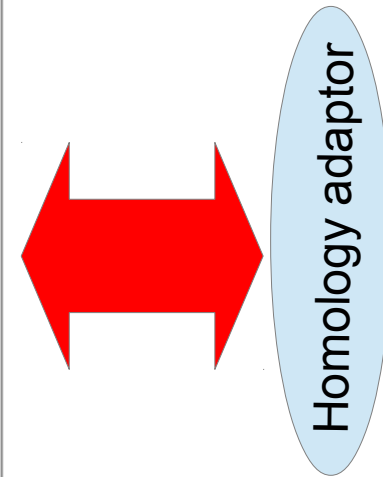
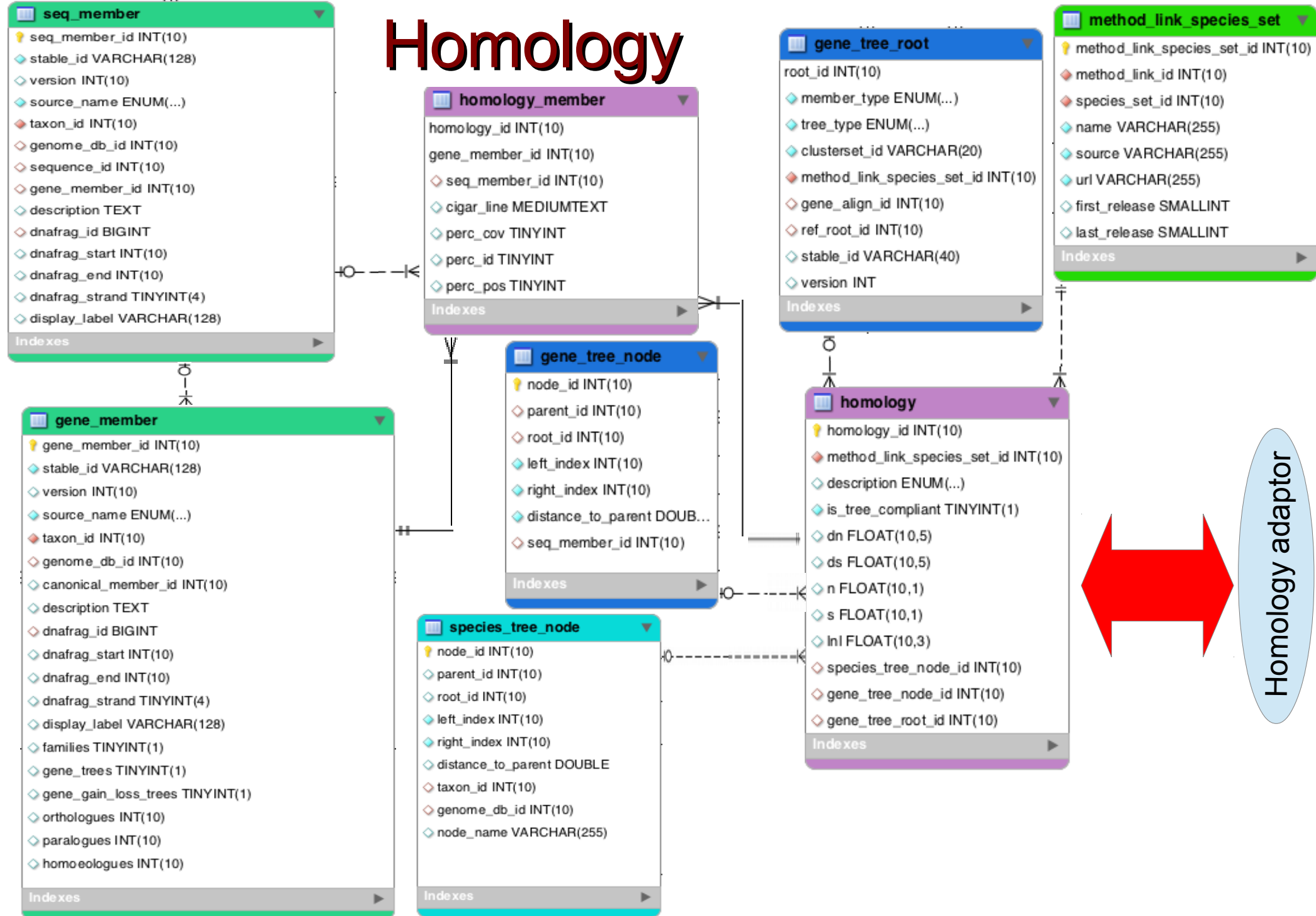
- Written in Object-Oriented Perl
- Used to retrieve data from and store data into the Ensembl Compara database
- Generalized to extend to non-Ensembl genomic data (Uniprot)
- Follows same 'Object Adaptor' & 'Data Object' design as the Core API



# Ensembl API Architecture



# Homology



# Compara template script

```
use strict;
use Bio::EnsEMBL::Registry;
my $reg = "Bio::EnsEMBL::Registry";

# Auto-configure the registry
$reg->load_registry_from_db(
    -host => "ensembl.ensembl.org",
    -user => "anonymous"
);

# Get the adaptor object for the data type you want
# e.g. GeneTree
my $xx_adaptor = $reg->get_adaptor("Multi", "compara", "XX");

# Fetch the data objects using the adaptor
# e.g. get all the genes in a given gene tree
my $all_interesting_xx = $xx_adaptor->fetch_all_by_YY();

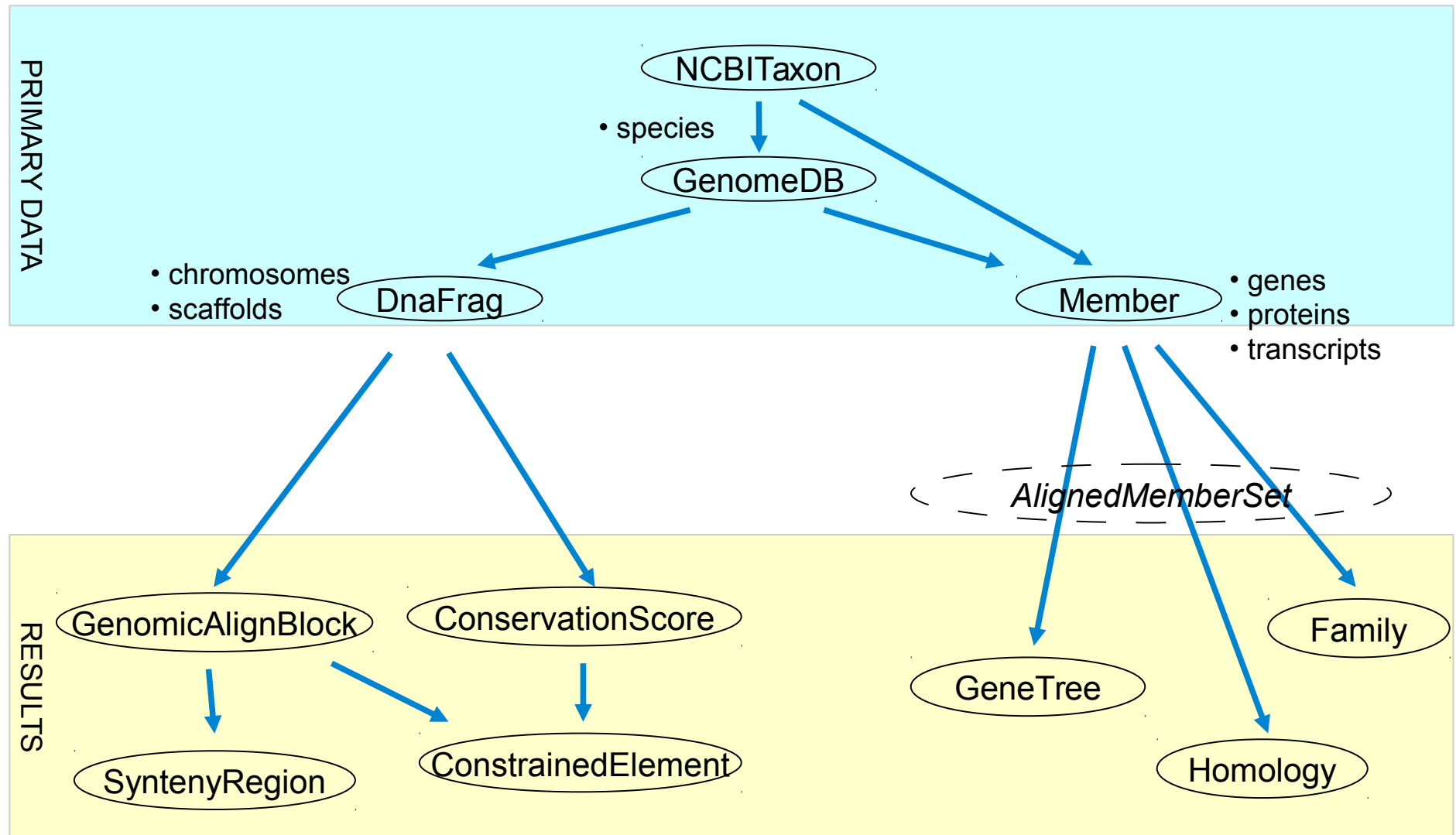
print "All XX objects from E!Compara :\n";
foreach my $this_xx (@$all_interesting_xx) {
    # Do some stuff with the data object
    print "\t", $this_xx->stable_id, "\n";
}
```

# Help & Useful documentation

- perldoc – Viewer for offline API documentation
  - `shell> perldoc Bio::Ensembl::Compara::GenomeDB`
  - `shell> perldoc Bio::Ensembl::Compara::DBSQL::DnaFragAdaptor`
- Online documents (website)
  - <http://e81.ensembl.org/info/docs/Doxygen/compara-api/index.html>
  - <http://e81.ensembl.org/info/docs/api/compara/index.html>
- Mailing lists:
  - [dev@ensembl.org](mailto:dev@ensembl.org)
  - [helpdesk@ensembl.org](mailto:helpdesk@ensembl.org)



# Compara object model overview



# Outline of the course

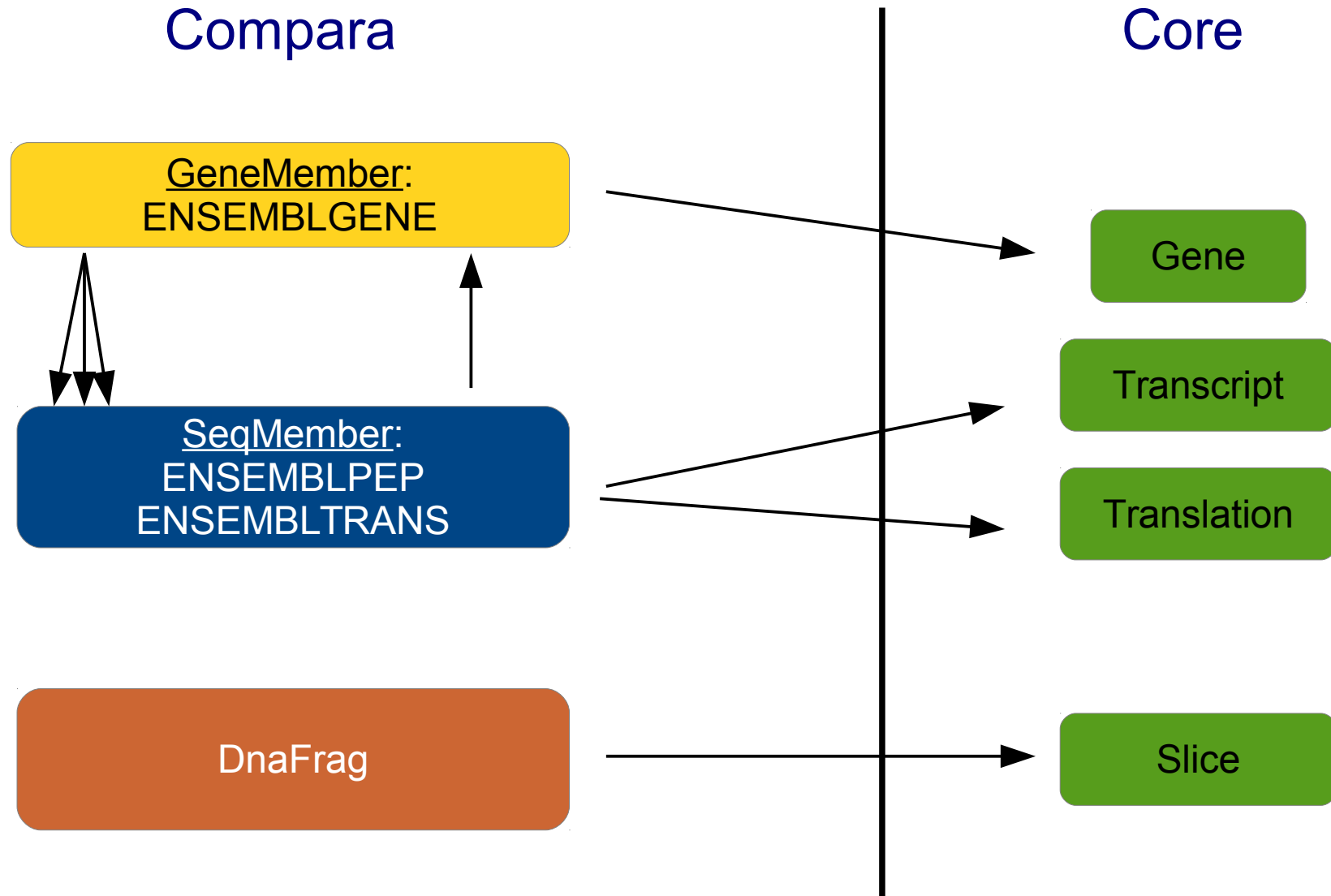
- Introduction about Compara
  - Resources
  - API
- Inputs
  - Species, Chromosomes, Genes
- Outputs
  - Gene analyses
  - Genome analyses



# Links between Compara and Core

- Compara only stores references to the Core objects
- The full data lies in the core databases

# Links between Compara and Core



# GenomeDB

- Represents a species  
Links the Compara database to the Core species databases

| Attributes  | Methods                                 |
|---|---|
| Species name  | <code>\$genomedb-&gt;name()</code>      |
| Assembly  | <code>\$genomedb-&gt;assembly()</code>  |
| Gene build  | <code>\$genomedb-&gt;genebuild()</code> |
| Taxon   | <code>\$genomedb-&gt;taxon_id()</code>  |
| <b>Adaptor methods</b>  |   |
| <code>\$genomedb_adaptor-&gt;fetch_all(...)</code>              |   |
| <code>\$genomedb_adaptor-&gt;fetch_by_registry_name(...)</code> |   |

# DnaFrag

- Represents a top-level region in the Compara database.
- Equivalent to a whole-sequence Slice

| Attributes   | Methods  |
|--|--|
| Region name  | <code>\$dnafrag-&gt;name()</code>              |
| Region type  | <code>\$dnafrag-&gt;coord_system_name()</code> |
| <b>Adaptor methods</b>   |  |
| <code>\$dnafrag_adaptor-&gt;fetch_all_by_GenomeDB_region(...)</code> |  |
| <code>\$dnafrag_adaptor-&gt;fetch_by_Slice(...)</code>               |  |
| <code>\$dnafrag_adaptor-&gt;fetch_by_GenomeDB_and_name(...)</code>   |  |

# Code Examples

## 1. GenomeDB

```
my $genome_db_adaptor = Bio::EnsEMBL::Registry->get_adaptor( "Multi", "compara", "GenomeDB");
my $list_ref_of_gdbs = $genome_db_adaptor->fetch_all();

foreach my $genome_db( @{ $list_ref_of_gdbs } ){
    print join( "\t", $genome_db->dbID(), $genome_db->name(), $genome_db->assembly() ), "\n";
}
```

## 2. DnaFrag

```
my $dnafrag_adaptor = $reg->get_adaptor("Multi", "compara", "DnaFrag");
my $gorilla_chr_dnafrags = $dnafrag_adaptor-
>fetch_all_by_GenomeDB_region( $gorilla_genome_db, 'chromosome' );

foreach my $dnafrag ( @{ $gorilla_chr_dnafrags } ){
    print "Chromosome ", $dnafrag->name(), " contains ", $dnafrag->length(), " bp.\n";
}
```

# Exercises – GenomeDB & DnaFrag

- Print the name, assembly version and genebuild version for all the GenomeDBs in the compara database
- Print all the chromosomes (DnaFrag) for chimpanzee



# GeneMember and SeqMember

- GeneMember for genes
  - `source_name`: ENSEMBLGENE
- SeqMember for RNAs and proteins
  - `source_name`: ENSEMBLPEP, ENSEMBLTRANS, Uniprot/SPTREMBL, Uniprot/SWISSPROT

| Attributes  | Methods   |
|---|---|
| Stable ID   | <code>\$member-&gt;stable_id()</code>   |
| Coordinates   | <code>\$member-&gt;chr_name()</code><br><code>\$member-&gt;chr_start() ...</code> |
| Sequence (SeqMember only)   | <code>\$member-&gt;sequence()</code>  |
| Function  | <code>\$member-&gt;description()</code>   |
| <b>Adaptor methods</b>  |   |
| <code>\$seq_member_adaptor-&gt;fetch_by_stable_id(...)</code>     |   |
| <code>\$gene_member_adaptor-&gt;fetch_all_by_GenomeDB(...)</code> |   |

# HOWTO: get an Ensembl ID from a gene symbol

- Compare only references genes by their Ensembl stable ID
- From a gene symbol, you first have to use the core API to get the stable id(s)
- Gene symbols may not be unique (for instance: U6)

```
# Get the Human gene adaptor
my $hg_adaptor = $reg->get_adaptor("human", "core", "Gene");

# Get all the genes
my $all_genes = $hg_adaptor->fetch_all_by_external_name(XX);

# For each gene
foreach my $gene (@{$all_genes}) {
    do some stuff with $gene->stable_id();
}
```

# Code Example - Member

```
my $seq_member_adaptor = $reg->get_adaptor("Multi", "compara", "SeqMember");
my $human_seq_members =
    $seq_member_adaptor->fetch_all_by_GenomeDB($gorilla_genome_db);

# print 10 peptide members and 10 transcript members
my ($pep_count, $trans_count) = (0, 0);
foreach my $seq_mem ( @{ $human_seq_members } ) {
    my $type = $seq_mem->source_name();
    if ( $type =~ m/PEP/ && $pep_count < 10 ){
        print $seq_mem->stable_id(), ":", $seq_mem->source_name(), "\n";
        $pep_count++;
    }
    elsif ( $type =~ m/TRANS/ && $trans_count < 10 ){
        print $seq_mem->stable_id(), ":", $seq_mem->source_name(), "\n";
        $trans_count++;
    }
    elsif ( $pep_count >= 10 && $trans_count >= 10 ) {
        last;
    }
}
```

# Exercises - *Member*

- Print the sequence of the Member corresponding to SwissProt protein O93279
- Find and print the sequence of all the peptide Members corresponding to the human protein-coding gene(s) FRAS1

# Outline of the course

- Introduction about Compara
  - Resources
  - API
- Inputs
  - Species, Chromosomes, Genes
- Outputs
  - Gene analyses
  - Genome analyses



# *AlignedMemberSet* object

- Base object that represents a set of members aligned together, e.g. a multiple alignment of peptides / ncRNAs
- “Applied” in gene trees, families, and homologies
- No specific adaptor

| Attributes                 | Methods  |
|----------------------------|--|
| List of members            | <code>\$aln-&gt;get_all_Members()</code><br><code>\$aln-&gt;get_all_GeneMembers()</code> |
| Alignment (BioPerl object) | <code>\$aln-&gt;get_SimpleAlign()</code>   |
| Description (if available) | <code>\$aln-&gt;description()</code>   |
| Stable ID (if available)   | <code>\$aln-&gt;stable_id()</code>   |

# HOWTO: print a BioPerl alignment

- Compara objects return alignments as BioPerl instances

```
$aln->get_SimpleAlign()
```

- BioPerl provides an AlignIO object to format the actual output in various formats (fasta, clustalw, phylip ... )

```
use Bio::AlignIO;
```

```
# Get the alignIO object from BioPerl
```

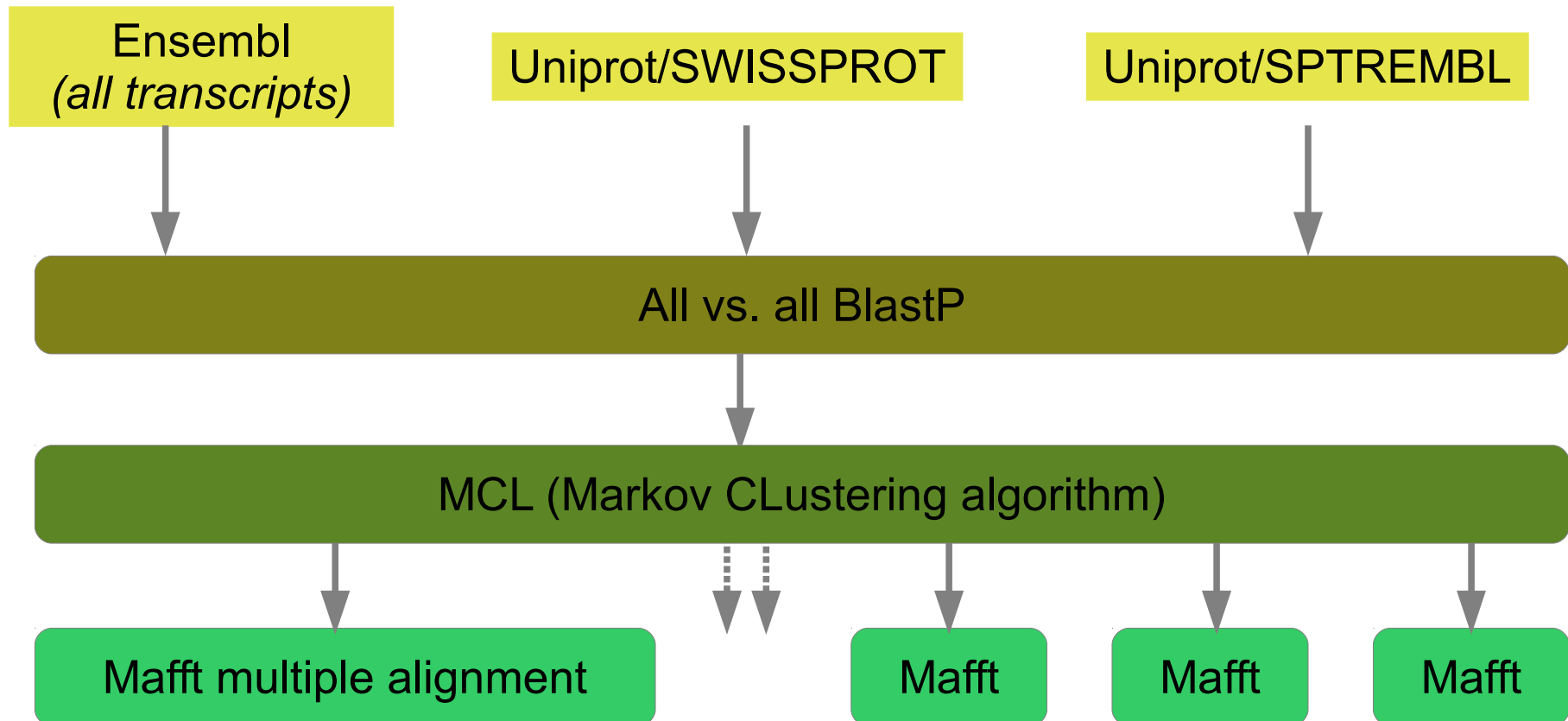
```
my $alignIO = Bio::AlignIO->newFh(-format => "fasta");
```

```
# Print the alignment
```

```
print $alignIO $aln;
```

# Families

Families are clusters of similar peptides



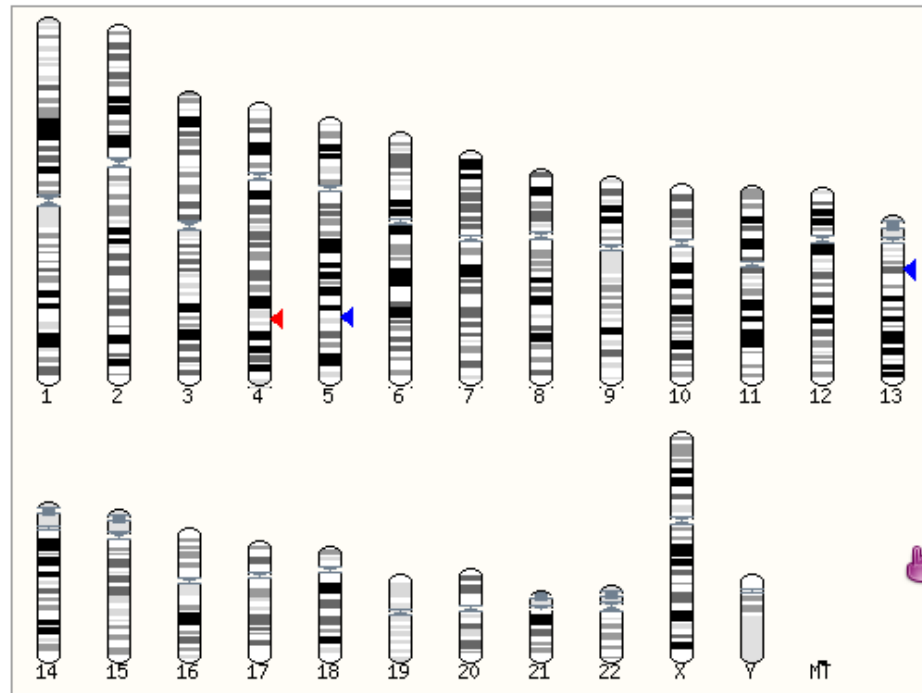


# Example on the web:

## ENSM00750001632338 in Human

HUMAN genes in this family

Ensembl genes containing proteins in family ENSM00750001632338



| Gene ID and Location                                     | Gene Name             | Description (if known)                             |
|--|-----------------------|--|
| <a href="#">ENSG00000170365</a><br>Chromosome 4: 146.40m | <a href="#">SMAD1</a> | SMAD family member 1 [Source:HGNC Symbol;Acc:6767] |
| <a href="#">ENSG00000113658</a><br>Chromosome 5: 135.47m | <a href="#">SMAD5</a> | SMAD family member 5 [Source:HGNC Symbol;Acc:6771] |
| <a href="#">ENSG00000120693</a><br>Chromosome 13: 37.42m | <a href="#">SMAD9</a> | SMAD family member 9 [Source:HGNC Symbol;Acc:6774] |

# Family object

- (almost) the same methods as in *AlignedMemberSet*
- Alternative transcripts can belong to different families !



| Attributes   | Methods                                     |
|--|---|
| Alignment  | <code>\$family-&gt;get_SimpleAlign()</code> |
| Biological function  | <code>\$family-&gt;description()</code>     |
| Gene content   | <code>\$family-&gt;get_all_Members()</code> |
| <b>Adaptor methods</b>   |   |
| <code>\$family_adaptor-&gt;fetch_all_by_GeneMember(...)</code> |   |
| <code>\$family_adaptor-&gt;fetch_by_SeqMember(...)</code>      |   |
| <code>\$family_adaptor-&gt;fetch_by_stable_id(...)</code>      |   |

# Code Example - Family

```
my $family_adaptor = $reg->get_adaptor("Multi", "compara", "Family");
my $ddx_families = $family_adaptor-
>fetch_by_description_with_wildcards('dead box', 1);

# print first 10 family descriptions
my $c = 0;
foreach my $fam ( @{ $ddx_families } ) {
    print $fam->description(), "\n";
    $c++;
    last if $c >= 10;
}
```

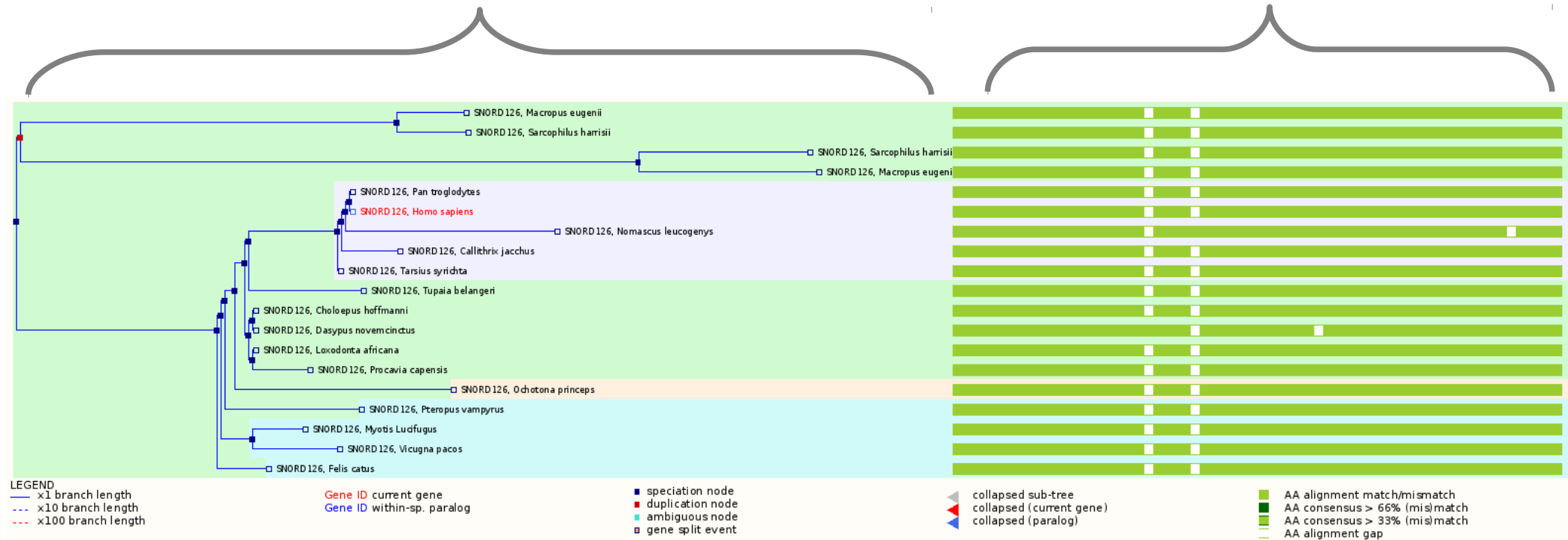
# Exercises - Families

- Get the multiple alignment corresponding to the family with the stable id ENSFM00250000006121
- Get the families predicted for the human gene ENSG00000139618. What do you notice ?

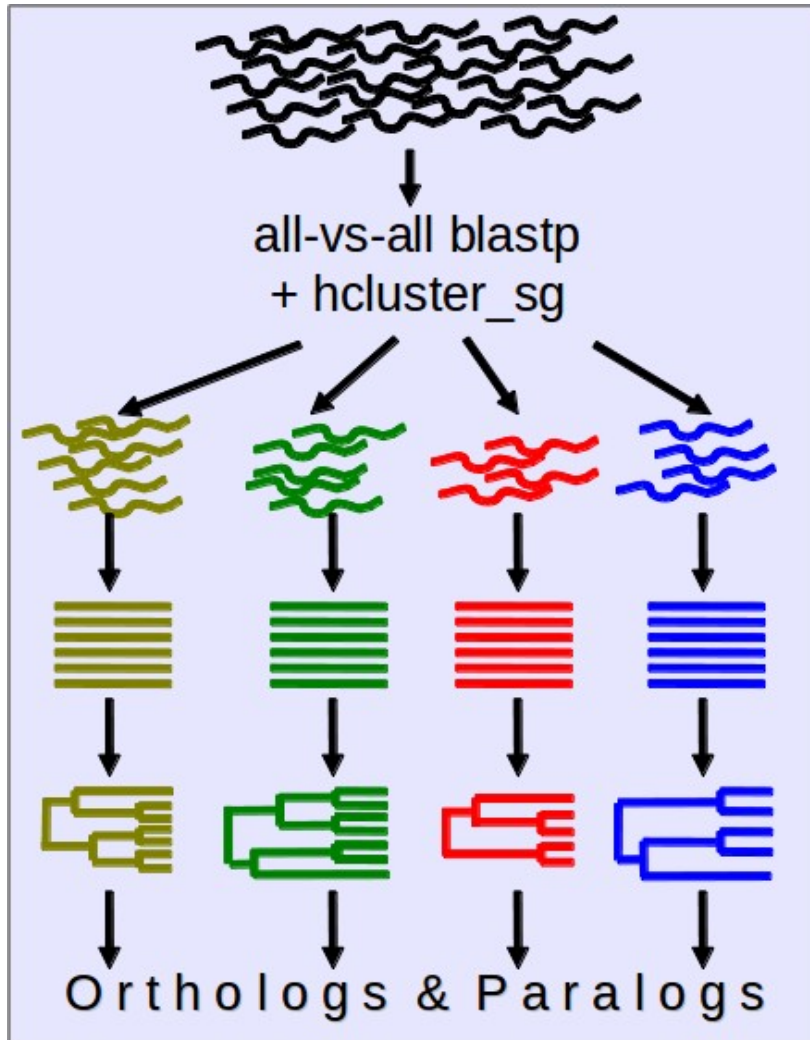
# GeneTree example on the website

## Tree

## Multiple alignment



# Protein-Tree pipeline overview



All *e!* genes – canonical prot.

BLAST

hcluster\_sg

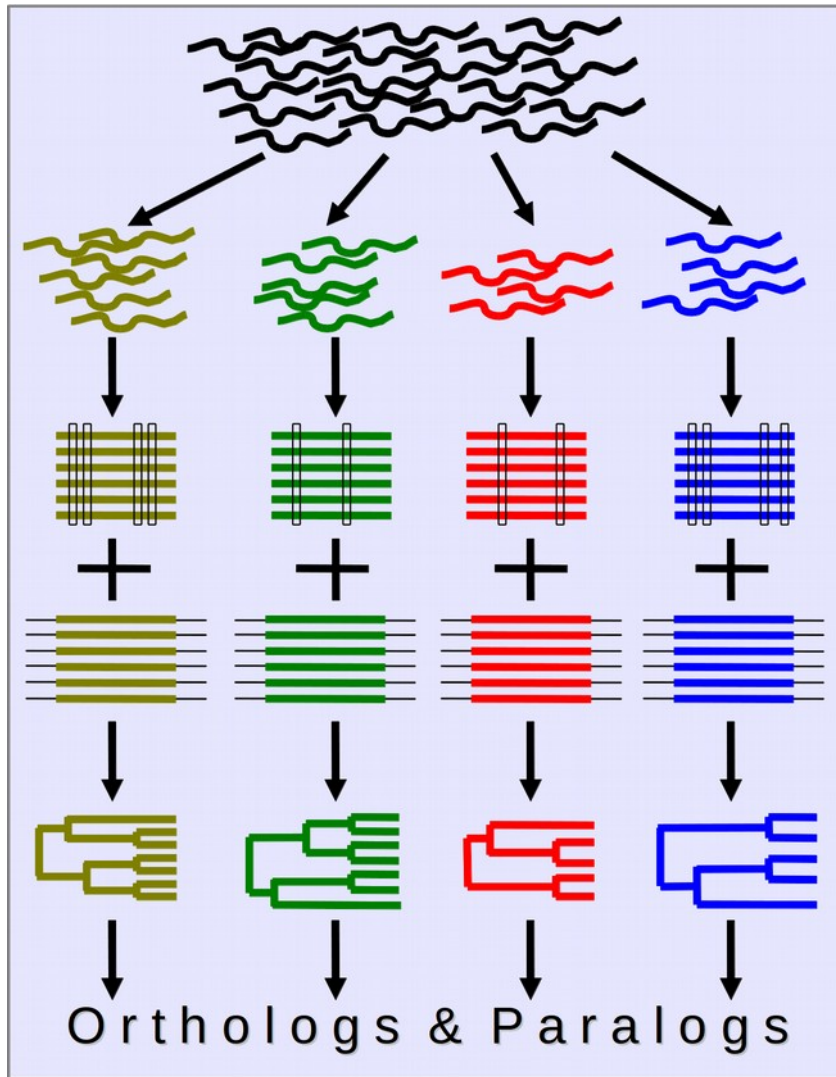
MCoffee: MSA

TreeBeST: (+ reconciliation)

Ortholog/Paralog inference

Vilella et al., Genome Res. 2009

# ncRNA-Tree pipeline overview



All *e!* ncRNA genes

Grouped in Family Models - RFAM

Infernal alignment + RaxML trees

PRANK alignment + NJ/ML trees

TreeBeST (tree reconciliation)

Ortholog/Paralog inference

Pignatelli et al., in preparation

# GeneTree object

- fetch\_all\** methods may require some more arguments:

```
-clusterset_id => 'default'  
-tree_type => 'tree'  
-member_type => 'protein' or 'ncrna'
```



| Attributes  | Methods   |
|---|---|
| Alignment   | <code>\$family-&gt;get_SimpleAlign()</code>   |
| Tree export   | <code>\$tree-&gt;newick_format('simple')</code><br><code>\$tree-&gt;nhx_format('full')</code><br><code>\$tree-&gt;print_tree()</code> |
| Stable ID   | <code>\$tree-&gt;stable_id()</code>   |
| <b>Adaptor methods</b>  |   |
| <code>\$genetree_adaptor-&gt;fetch_by_stable_id(...)</code>       |   |
| <code>\$genetree_adaptor-&gt;fetch_default_for_Member(...)</code> |   |

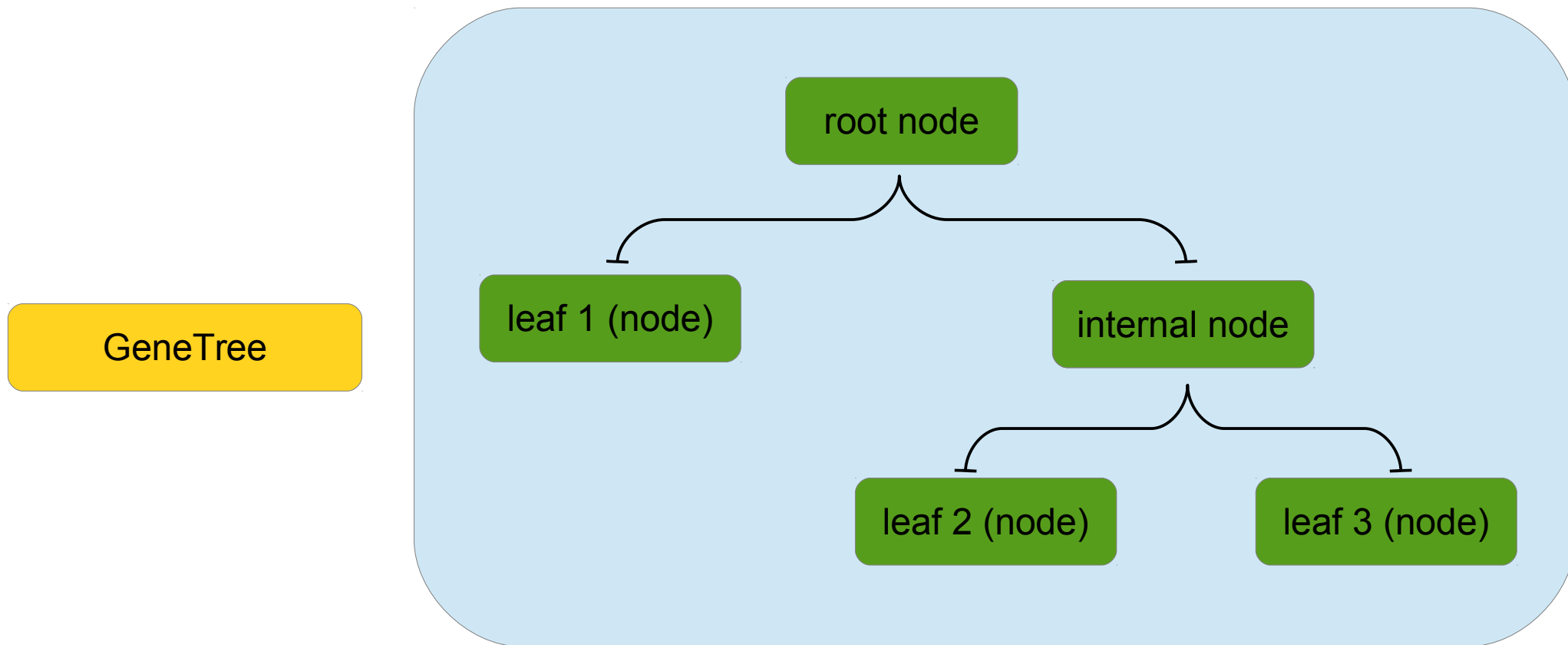


# Exercises – Protein and ncRNA trees

- Print the protein tree with the stable id ENSGT00390000003602
- Print all the members of the tree containing the human ncRNA gene ENSG00000238344, and their alignment

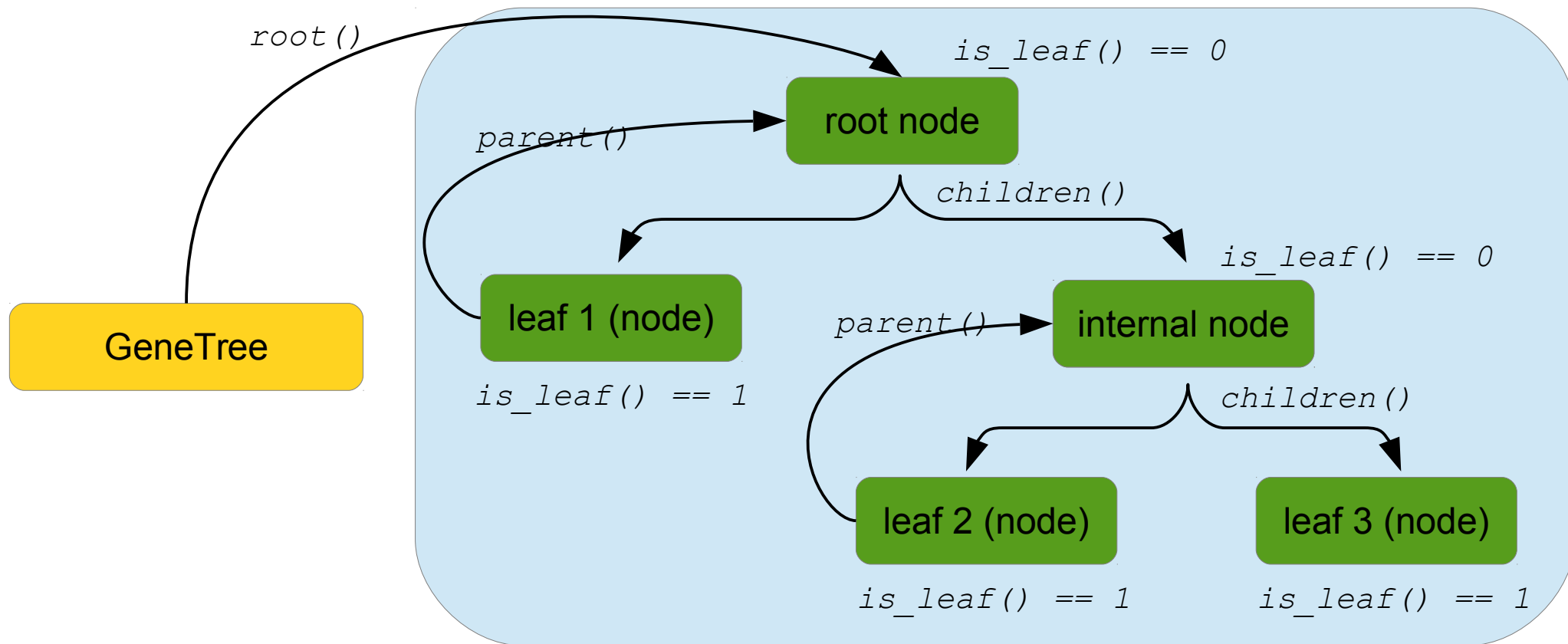
# GeneTreeNode object

The actual tree structure is a hierarchy of *GeneTreeNode* objects



# GeneTreeNode object

The actual tree structure is a hierarchy of *GeneTreeNode* objects



## Extra information

`$node->node_type()`

`$node->duplication_confidence_score()`

`$node->taxonomy_level()`

`$node->bootstrap()`

# Outline of the course

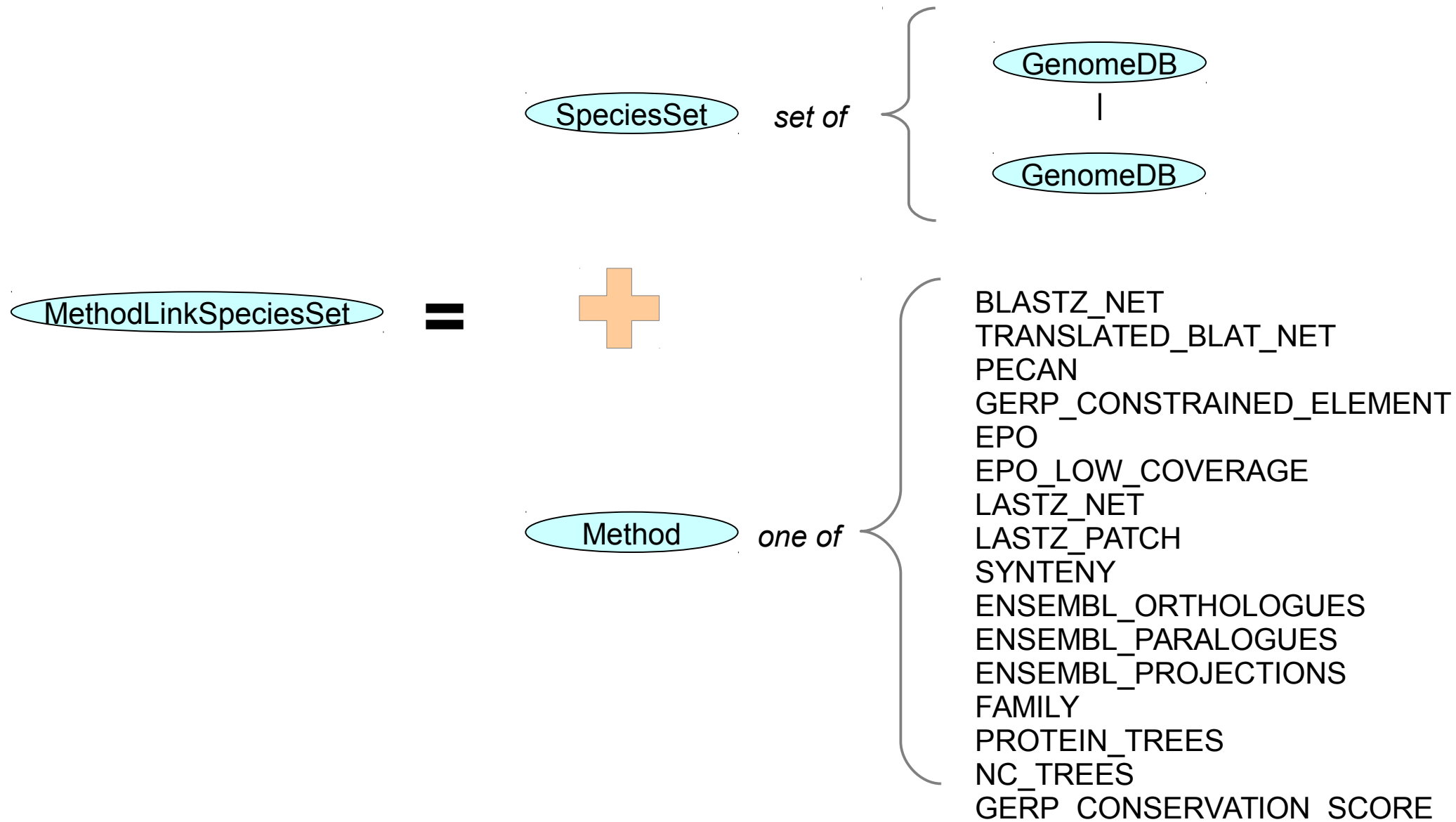
- Introduction about Compara
  - Resources
  - API
- Inputs
  - Species, Chromosomes, Genes
- Outputs
  - Gene analyses
  - Genome analyses



# MethodLinkSpeciesSet object

- The Compara database contains **lots** of cross-species comparisons
- There are multiple comparisons of the same type (pairwise alignments, homologies, etc)
- We need a way of defining which analysis is performed on which genomes
- Many adaptor methods require a MethodLinkSpeciesSet

# MethodLinkSpeciesSet object



# MethodLinkSpeciesSet

- Links a method (an analysis) to a set of species

| Attributes   | Methods                                     |
|--|---|
| Name   | <code>\$mlss-&gt;name()</code>              |
| Type of analysis   | <code>\$mlss-&gt;method()-&gt;type()</code> |
| List of GenomeDBs  | <code>\$mlss-&gt;species_set()</code>       |
| <b>Adaptor methods</b>   |   |
| <code>\$mlss_adaptor-&gt;fetch_by_method_link_type_registry_aliases()</code> |   |
| <code>\$mlss_adaptor-&gt;fetch_by_method_link_type_species_set_name()</code> |   |

# Example Code – MethodLinkSpeciesSet

```
my $gdb_a = $reg->get_adaptor( "Multi", "compara", "GenomeDB" );
my $gorilla_genome_db = $gdb_a->fetch_by_dbID(123);

my $mlss_adaptor = $reg->get_adaptor("Multi", "compara", "MethodLinkSpeciesSet");
my $gorilla_mlss_list = $mlss_adaptor->fetch_all_by_GenomeDB( $gorilla_genome_db );

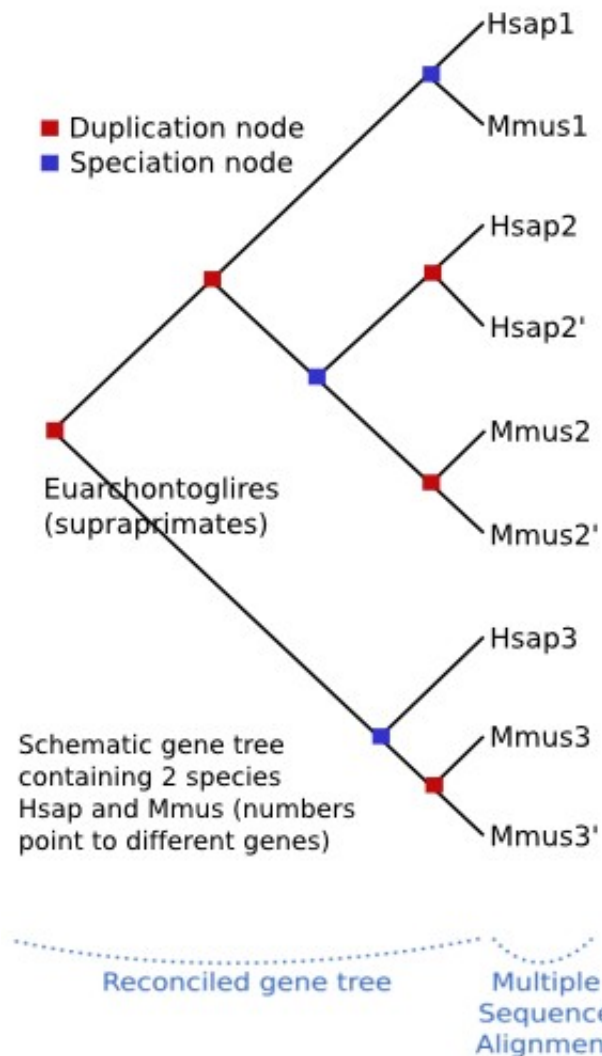
my $c = 0;
foreach my $mlss ( @{ $gorilla_mlss_list } ) {
    print join( "\t", $mlss->dbID(), $mlss->method->type() ), "\n";
    $c++;
    last if $c >= 10;
}
```



# Exercises – MethodLinkSpeciesSet

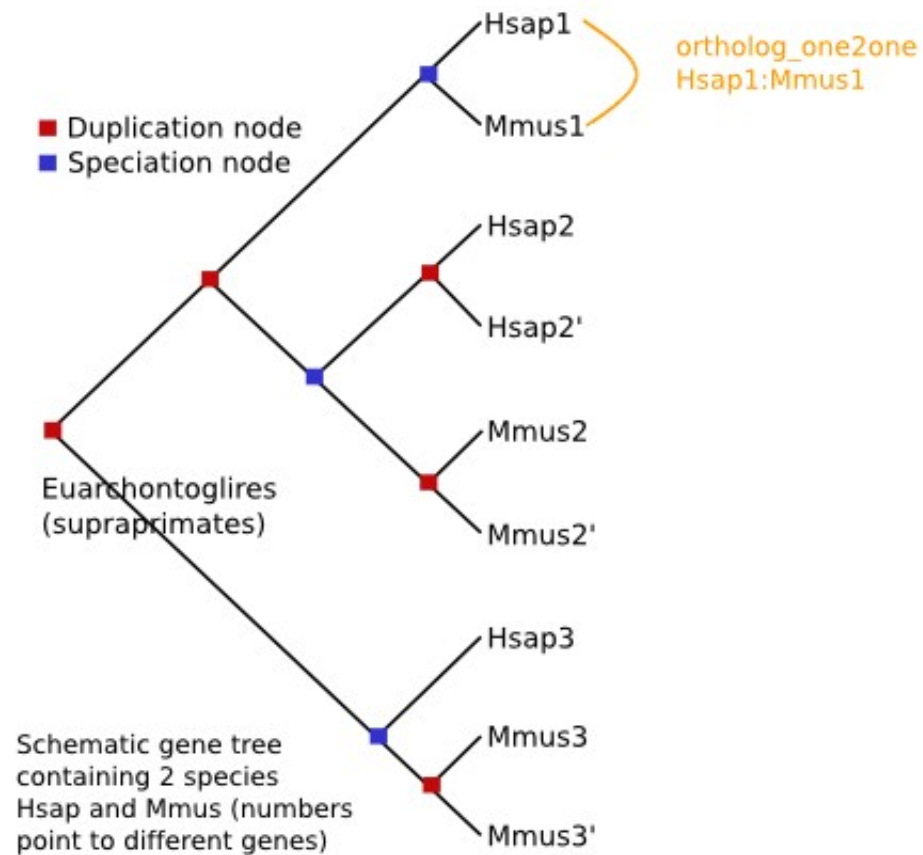
- Print the total number of MethodLinkSpeciesSet entries stored in the database
- Print a unique list of method\_link\_types and a count of their number in the database.
- Print the list of the species for the 17 eutherian mammals EPO alignments

# Homology inference

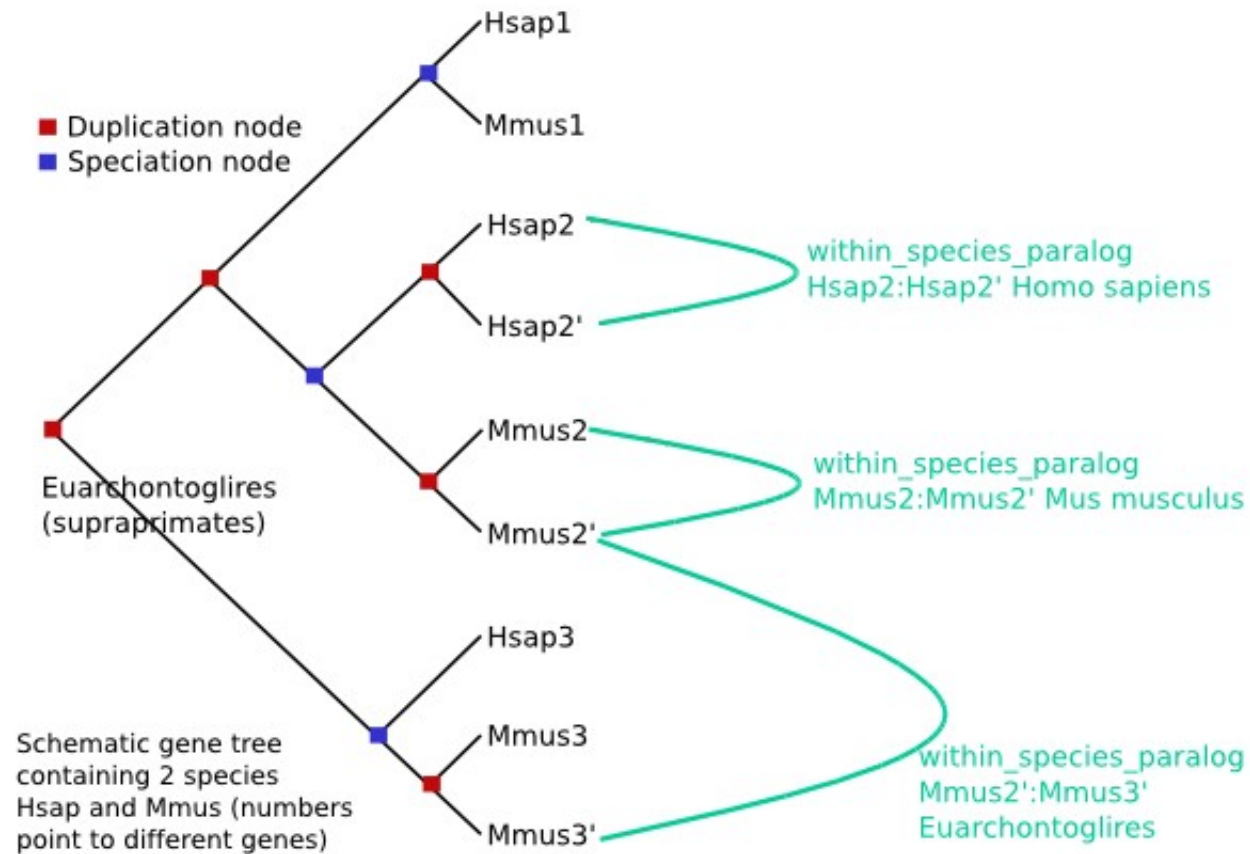


Consists in tagging the pairs of genes of all the trees with a relation type, depending on the tree topology.

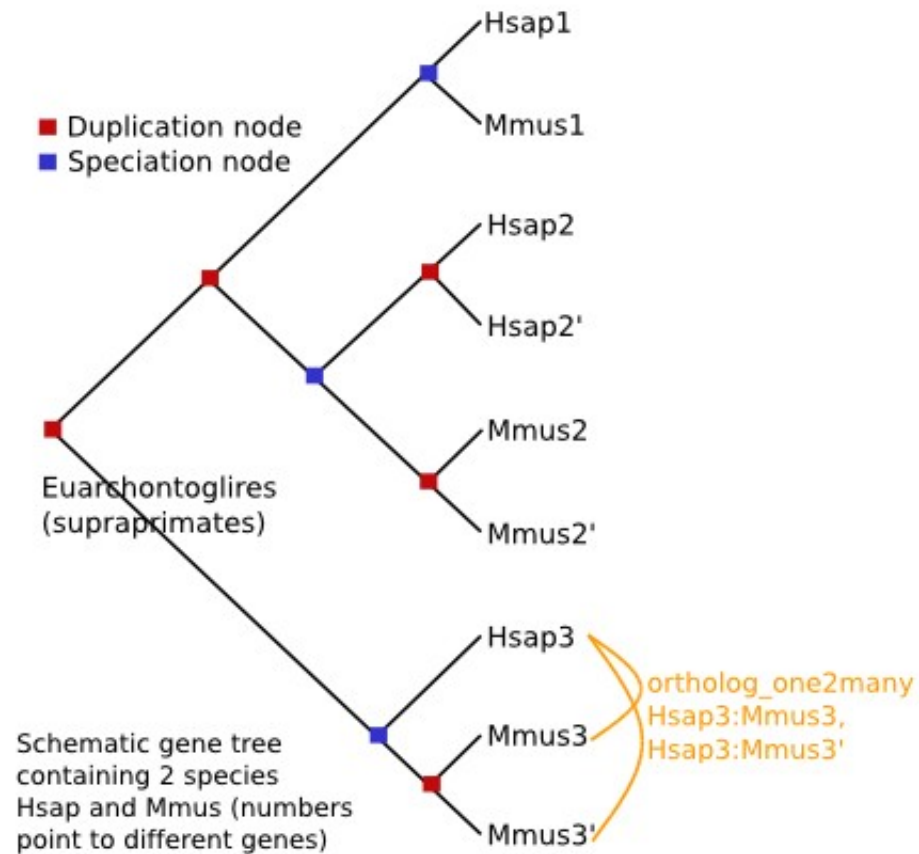
# Homology inference



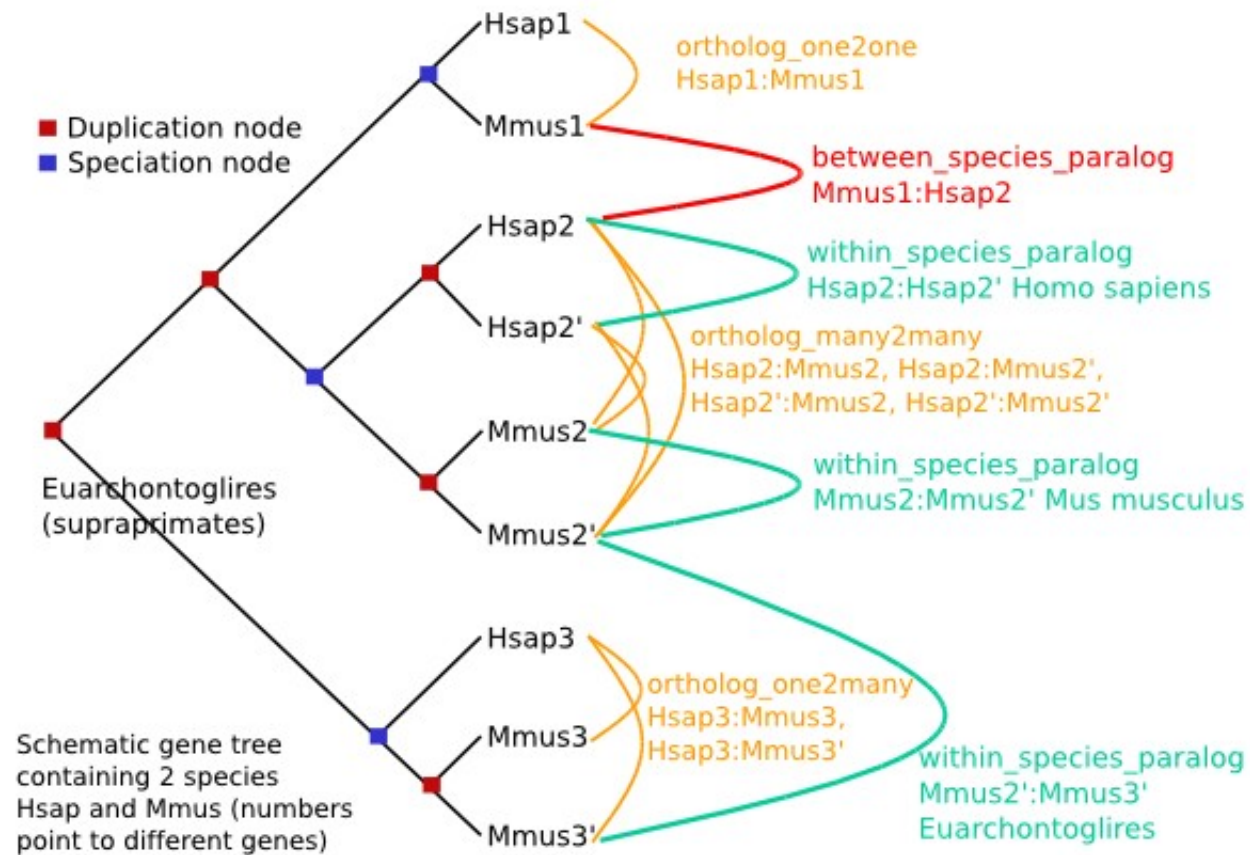
# Homology inference



# Homology inference



# Homology inference



# Homology object

- An Homology object links two genes together
- One-to-many relationships are split:
  - “H ortholog to M1” and “H ortholog to M2” are different objects



| Attributes  | Methods   |
|---|---|
| Alignment   | <code>\$homology-&gt;get_SimpleAlign()</code>   |
| Natural selection   | <code>\$homology-&gt;dn()</code> / <code>\$homology-&gt;ds()</code>                       |
| Gene content  | <code>\$homology-&gt;get_all_GeneMembers()</code>   |
| Homology characteristics  | <code>\$homology-&gt;description()</code><br><code>\$homology-&gt;taxonomy_level()</code> |
| <b>Adaptor methods</b>  |   |
| <code>\$homology_adaptor-&gt;fetch_all_by_Member(...)</code>                |   |
| <code>\$homology_adaptor-&gt;fetch_all_by_MethodLinkSpeciesSet(...)</code>  |   |
| <code>\$homology_adaptor-&gt;fetch_all_by_Member_paired_species(...)</code> |   |

# Code Example - Homology

```
my ( $human_genomedb_id, $gorilla_genomedb_id ) = ( 150, 123 );

my $mlss_adaptor = $reg->get_adaptor("Multi", "compara", "MethodLinkSpeciesSet");
my $homology_adaptor = $reg->get_adaptor("Multi", "compara", "Homology");

my $mlss = $mlss_adaptor->fetch_by_method_link_type_genome_db_ids(
    'ENSEMBL_ORTHOLOGUES',
    [$human_genomedb_id, $gorilla_genomedb_id]
);
my $orthologs = $homology_adaptor->fetch_all_by_MethodLinkSpeciesSet($mlss);

my $c = 0;
foreach my $orth ( @{ $orthologs } ){
    print $orth->toString(), "\n" if $orth->is_tree_compliant();
    $c++;
    last if $c >= 10;
}
```



# Exercises – Homologies (and MethodLinkSpeciesSet)

- Get all the homologues for the human gene ENSG00000229314
- Count the number of “one2one” homologues between human and mouse
- Find the human orthologues of ENSMUSG000000004843 and ENSMUSG000000025746. For each homology, display the alignment and the dn value. Comment on the divergence

# Outline of the course

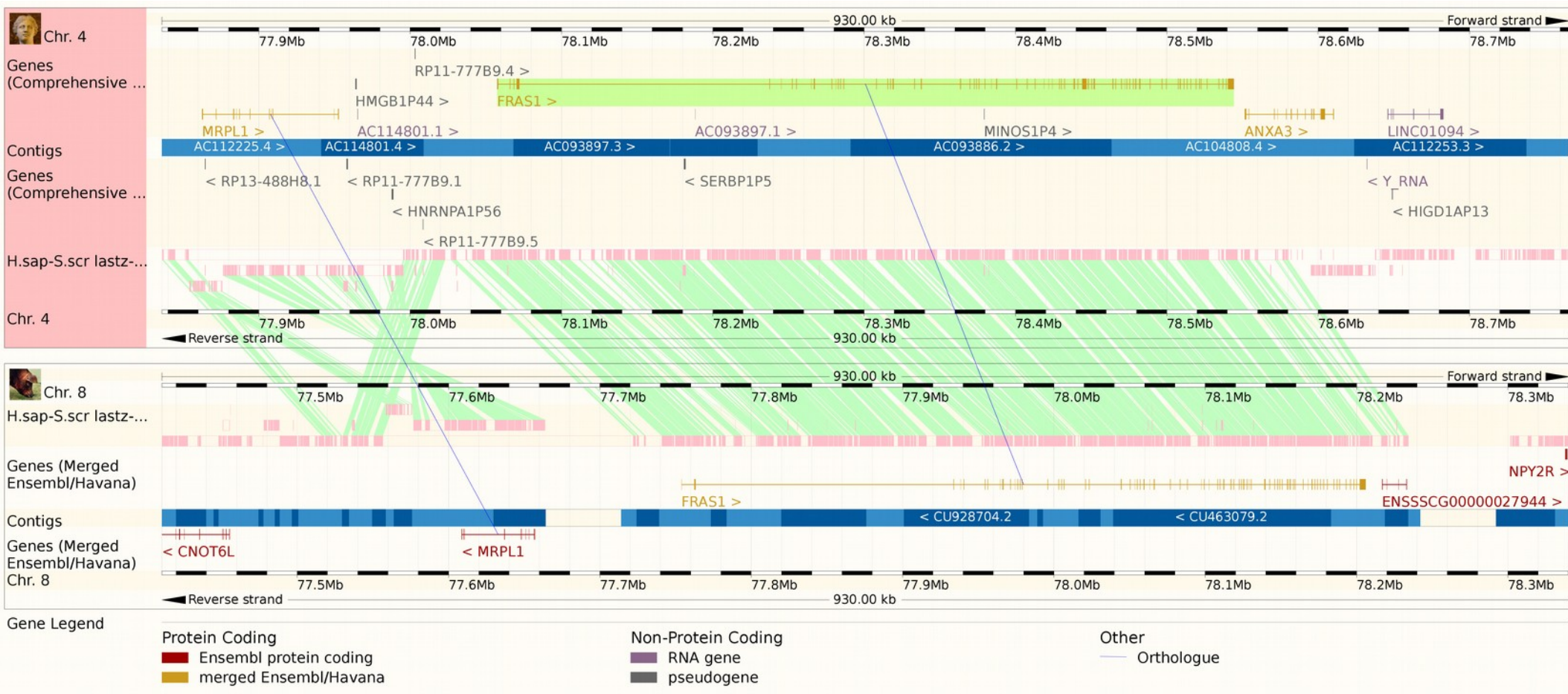
- Introduction about Compara
  - Resources
  - API
- Inputs
  - Species, Chromosomes, Genes
- Outputs
  - Gene analyses
  - Genome analyses



# Whole-genome alignments

## Alignments at the DNA level

Example: Human vs pig



# How are alignments stored ?

A small example :

|                             |                              |
|-----------------------------|------------------------------|
| gorilla_gorilla/MT/935-953  | gacat-ttaactaaaac-ccc        |
| macaca_mulatta/MT/1469-1488 | aacatcttaactaaacg-ccc        |
| pan_troglodytes/MT/934-953  | gatac-ttaacttaaaccccc        |
| pongo_pygmaeus/MT/940-958   | actac-ctaactaaaac-ccc        |
| homo_sapiens/MT/1516-1534   | gacat-ttaactaaaac-ccc        |
|                             | *       *****   **       *** |

GACATTTAACTAAAACCCC  
AACATCTTAACTAAACGCCC  
GATACTTAACTTAAACCCCC  
ACTACCTAACTAAAACCCC  
GACATTTAACTAAAACCCC

Sequences from core

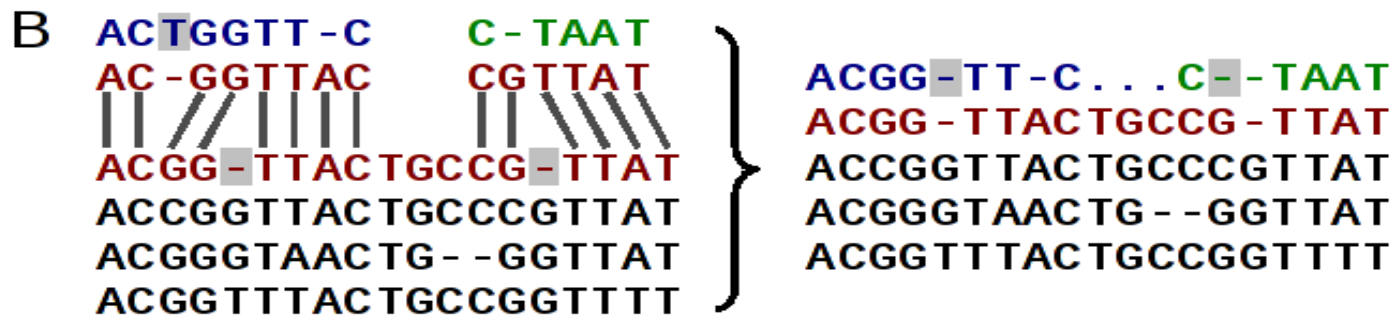
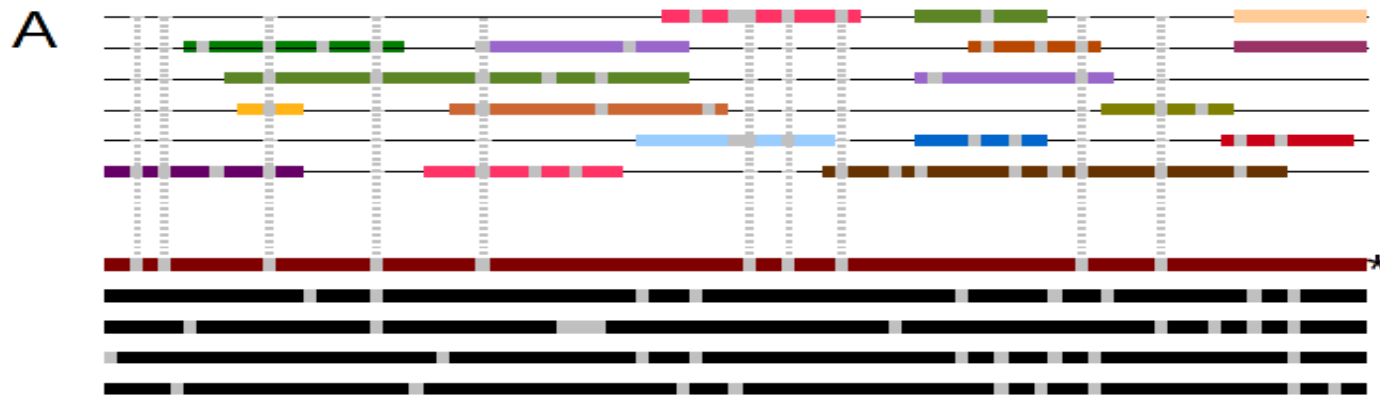
5MD11MD3M  
17MD3M  
5MD15M  
5MD11MD3M  
5MD11MD3M

CIGAR line in Compara

5 *genomic\_align* entries  
1 *genomic\_align\_block* entry

# Adding low-coverage genomes

- Low coverage genomes cannot be fully assembled
- Resulting assembly is too scattered to be used with Enredo
- Run EPO on high-coverage genomes only
- Map 2X genomes using pairwise alignments on a reference species



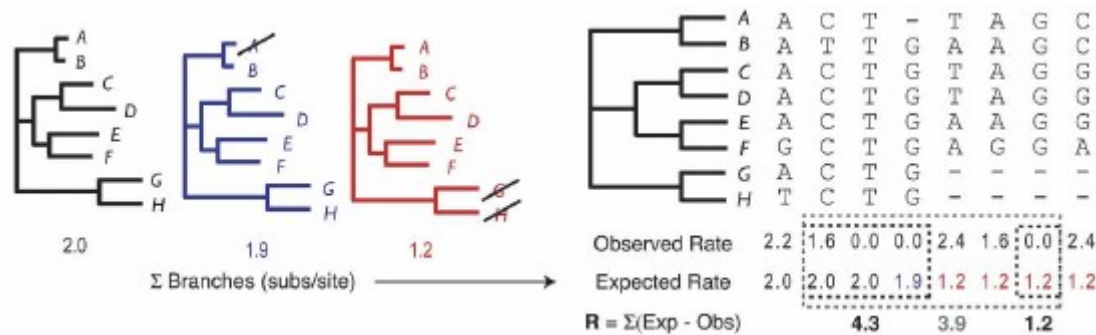
# Objects on the genomic side

- A GenomicAlignBlock represents an alignment between two or more regions of genomic DNA. Within these blocks every region of genomic DNA is represented by a GenomicAlign object.
- A ConstrainedElement represent regions in the multiple alignment which appear to be under functional constraint.
- Synteny blocks are derived from Lastz-net alignments
  - group syntenic alignments closer than 200Kb
  - link syntenic groups closer than 3Mb
  - minimum length of the syntenic block: 100 kb



# GERP Constrained Elements

- Stretches of the alignment with a high conservation



Cooper et al. *Genome Research*, 2005

- Constrained elements and coding exons
  - 74% of coding exons are associated with constr. elem.
  - 22% of constr. elem. are associated with coding exons

# GenomicAlignBlock

- An alignment-block (across 2 or more sequences)
- The adaptor returns the blocks that **overlap** the query region  
→ Call `restrict_between_reference_positions()`

| Attributes   | Methods   |
|--|---|
| BioPerl alignment object   | <code>\$gab-&gt;get_SimpleAlign()</code>                      |
| Aligned sequences  | <code>\$gab-&gt;get_all_GenomicAligns()</code>                |
| (Restrict the block)   | <code>\$gab-&gt;restrict_between_reference_positions()</code> |
| Adaptor methods  |   |
| <code>\$gab_adaptor-&gt;fetch_all_by_MethodLinkSpeciesSet_Slice()</code> |   |

- GenomicAlign has a similar interface to Members, e.g.  
`$ga→dnafrag`, `$ga→dnafrag_start`, etc



# Code Example - GenomicAlignBlock

```
my $mlss_adaptor = $reg->get_adaptor("Multi", "compara", "MethodLinkSpeciesSet");
my $genomic_align_block_adaptor = $reg->get_adaptor("Multi", "compara",
"GenomicAlignBlock");

my $epo_mlss_list = $mlss_adaptor->fetch_all_by_method_link_type("EPO");

foreach my $epo_mlss ( @{ $epo_mlss_list } ){
    my $genomic_align_blocks = $genomic_align_block_adaptor->
        fetch_all_by_MethodLinkSpeciesSet( $epo_mlss );
    print $epo_mlss->name() . " has " . scalar( @{ $genomic_align_blocks } ) .
        " alignment blocks\n";
}
```

# Exercises – Genomic Alignments

- Print the LASTZ-NET alignments for pig chromosome 15 with cow (using pig coordinates 105734307 and 105739335).
- Change the above example so that it prints the 17-way eutherian mammal (EPO) multiple alignments.
- Print the constrained element alignments from the above pig locus (use the constrained elements generated from the EPO\_LOW\_COVERAGE mammals alignments)

# Exercises – Synteny

- Print the pig-cow synteny map using pig chromosome 15 as a reference

*web reference:*

[http://www.ensembl.org/Sus\\_scrofa/Location/Synteny?  
r=15&otherspecies=Bos\\_taurus](http://www.ensembl.org/Sus_scrofa/Location/Synteny?r=15&otherspecies=Bos_taurus)

# Acknowledgements



Leo Mateus Matthieu Carla Aj



D48–D55 *Nucleic Acids Research*, 2013, Vol. 41, Database issue  
doi:10.1093/nar/gks1236

Published online 30 November 2012

## Ensembl 2013

Paul Flicek<sup>1,2,\*</sup>, Ikhlak Ahmed<sup>1</sup>, M. Ridwan Amode<sup>2</sup>, Daniel Barrell<sup>2</sup>, Kathryn Beal<sup>1</sup>, Simon Brent<sup>2</sup>, Denise Carvalho-Silva<sup>1</sup>, Peter Clapham<sup>2</sup>, Guy Coates<sup>2</sup>, Susan Fairley<sup>2</sup>, Stephen Fitzgerald<sup>1</sup>, Laurent Gil<sup>1</sup>, Carlos García-Girón<sup>2</sup>, Leo Gordon<sup>1</sup>, Thibaut Hourlier<sup>2</sup>, Sarah Hunt<sup>1</sup>, Thomas Juettemann<sup>1</sup>, Andreas K. Kähäri<sup>2</sup>, Stephen Keenan<sup>1</sup>, Monika Komorowska<sup>1</sup>, Eugene Kulesha<sup>1</sup>, Ian Longden<sup>1</sup>, Thomas Maurel<sup>1</sup>, William M. McLaren<sup>1</sup>, Matthieu Muffato<sup>1</sup>, Rishi Nag<sup>2</sup>, Bert Overduin<sup>1</sup>, Miguel Pignatelli<sup>1</sup>, Bethan Pritchard<sup>2</sup>, Emily Pritchard<sup>1</sup>, Harpreet Singh Riat<sup>2</sup>, Graham R. S. Ritchie<sup>1</sup>, Magali Ruffier<sup>1</sup>, Michael Schuster<sup>1</sup>, Daniel Sheppard<sup>2</sup>, Daniel Sobral<sup>1</sup>, Kieron Taylor<sup>1</sup>, Anja Thormann<sup>1</sup>, Stephen Trevanion<sup>2</sup>, Simon White<sup>2</sup>, Steven P. Wilder<sup>1</sup>, Bronwen L. Aken<sup>2</sup>, Ewan Birney<sup>1</sup>, Fiona Cunningham<sup>1</sup>, Ian Dunham<sup>1</sup>, Jennifer Harrow<sup>2</sup>, Javier Herrero<sup>1</sup>, Tim J. P. Hubbard<sup>2</sup>, Nathan Johnson<sup>1</sup>, Rhoda Kinsella<sup>1</sup>, Anne Parker<sup>2</sup>, Giulietta Spudich<sup>1</sup>, Andy Yates<sup>1</sup>, Amonida Zadissa<sup>2</sup> and Stephen M. J. Searle<sup>2</sup>

<sup>1</sup>European Bioinformatics Institute, Wellcome Trust Genome Campus, Hinxton Cambridge CB10 1SD, UK and

<sup>2</sup>Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge CB10 1SA, UK

## Funding

wellcome trust



Co-funded by the  
European Union

