

# Bespoke comparative genomics software architectures a case-study for future HPC on Aristotelis

*Biosciences using High Performance Computing (HPC) Systems*



Christos A. Ouzounis [cao@csd.auth.gr](mailto:cao@csd.auth.gr) • 19 jun 2025





# HPC for Bioinformatics

cgg\_toolkit Hands-On Demo

need a short intro to concepts?  
and some elements of HPC

A screenshot of the CGG Computational Genomics Group Services website. The page features a navigation bar with links to Services, Projects, Publications, Documentation, Sponsors, People, Collaborators, and Internal. The main content area is titled 'Key software' and lists several tools: MagicMatch, BioLayout, geneCAST, geneRAGE, and TribeMCL, each with a brief description and a PubMed reference. Below this, there are sections for 'Key services 1: Complete Genome Sequences' and 'Key services 2: Genome Comparison using CoGenT'. The background of the website has a network diagram with yellow nodes and blue lines.

**CGG** Computational Genomics Group | **Services**

[Services](#) [Projects](#) [Publications](#) [Documentation](#) [Sponsors](#) [People](#) [Collaborators](#) [Internal](#)

### Key software

**MagicMatch**  
An efficient method to map sequence identifiers across databases [PubMed:15961438](#)

**BioLayout**  
An automatic graph layout algorithm for similarity visualization [PubMed:16000016](#)

**geneCAST**  
An algorithm for the complexity analysis of sequence tracts - filters and masks database query sequences [PubMed:11120681](#)

**geneRAGE**  
An algorithm for sequence clustering and automated domain detection [PubMed:10871267](#)

**TribeMCL**  
An efficient algorithm for large-scale detection of protein families [PubMed:11917018](#)

### Key services 1: Complete Genome Sequences

**CoGenT++**  
The CoGenT++ sitemap, clickable extended services [PubMed:16216832](#)

**CoGenT**  
The Complete GENome Tracking database [PubMed:12874064](#)

**GenMed**  
Continuous tracking of genomes in CoGenT by PubMed abstracts [PubMed:15864286](#)

**iCAST**  
Interactive detection and masking of low-complexity regions [PubMed:16216832](#)

**BlastServer**  
A BLAST service against CoGenT [PubMed:16216832](#)

### Key services 2: Genome Comparison using CoGenT



# *life is complex matter*

*an information universe*

- *biosciences: PBytes of information with large-scale installations*
- *analysis of these data => top-five players of Big Data*
- *along with financial, military, comms and physics/astronomy*
- *noisy, multi-scale: ideal playground for complex systems*

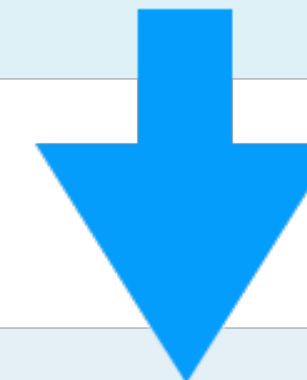




# central dogma from DNA to protein

- DNA contains 'genes'
- .. that encode
- .. (via an intermediate, mRNA)
- .. proteins
- both have distinct alphabets
  - DNA: 4 characters {A C G T}
  - proteins: 20 characters {aa's}

```
ATGCCCGTTGCCCACGTTGCCTTGCCCGTTCGCTTCCCTCGTACCTTTGACTATCTGCTGCCAGAAGGCATGACGGTTAAAGCTGGGTGTCGCGTGCGCGT
GCCGTTTGGCAAACAGCAGGAGCGCATCGGGATTGTGGTATCAGTTAGCGATGCCAGCGAACTGCCGCTCAATGAGCTAAAAGCGGTAGTCGAAGTGCTGG
ATAGTGAGCCGGTGTTTACTCACTCCGCTCTGGCGATTGCTGCTATGGGCGGCAGATTACTATCATCATCCGATTGGCGATGTGCTGTTTCATGCCCTTGCCG
ATTTTACTACGCCAGGGGCGGCCTGCGGCGAACGCGCCGATGTGGTACTGGTTTGCCACTGAACAAGGCCAGGCGGTGGATCTGAACAGCCTGAAACGCTC
CCCCAAGCAACAACAGGCGCTGGCGGCGTTACGGCAAGGCCAAAATCTGGCGCGACCAGGTCGCCACGCTCGAATTTAATGATGCCGCGTTGCAAGCGCTAC
GCAAAAAGGTCTGTGTGATTTAGCAAGTGAACACCAGAGTTTAGCGACTGGCGAACGAACATATGCCGTTTCTGGTGAGCGGTTGCGATTGAATACCGAA
CAGGCCACCGCCGTTGGCGCAATTCATAGCGCGGCAGATACTTTTTCTGCCTGGCTGCTGGCGGGCGTTACCGGTTCCGGTAAAACGGAGGTTTATCTCAG
CGTACTGGAAAACGTGCTCGCTCAGGGCAAACAGGCGCTGGTGATGGTGCCGAAATCGGCCTGACACCGCAAACCTATCGCCCGTTTTCTGTAACGTTTTA
ATGCCCCCGTGGAAGTTCTGCATTCCGGCCTGAACGACAGCGAGCGTCTTTTCGGCGTGGCTGAAAGCGAAAAATGGTGAGGCGGCGATTGTGATCGGCACC
CGCTCCGCGCTGTTTACGCCGTTTAAAAATCTCGGCGTGATTGTTCATTGATGAAGAGCACGACAGCTCCTACAAGCAGCAGGAAGGCTGGCGCTATCATGC
CCGCGACCTGGCGGTGTATCGTGCGCACAGCGAGCAAATCCCGATTATTCTTGGCTCCGCAACGCCCGCGCTGGAAACGTTATGCAACGTCCAGCAGAAAA
AATACCGCCTGCTGCGCCTGACCCGTCGGGCGAGGAATGCGCGTCCGGCAATTCACATGTGCTGGATTAAAAGGTGAGAAGGTGCAGGCAGGTCTGGCT
CCGGCGTTAATCACTCGTATGCGCCAGCATTACAGGCTGATAACCAGGTCACTCTCTTTCTTAACCGCCGTGGCTTTGCGCCTGCACTGCTGTGCCACGA
CTGTGGCTGGATTGCCGAATGCCACGTTGCGATCACTACTACACGCTGCATCAGGCGCAGCACCATCTGCGCTGCCACCACTGTGACAGTCAGCGTCCGG
TGCCGCGCCAGTGCCCTTCTGCGGTTCCACGCACCTGGTCCCCGTGGGGCTGGGCACCGAACAGCTTGAACAGACGCTCGCGCCGTTGTTCCCCGGCGTG
CCATTTCTCGTATCGACCGCGATACCACAGCCGCAAAGGGGCGCTGGAACAGCAACTGGCAGAAGTACATCGCGGCGGCGCGCGGATTTTGATTGGTAC
ACAAATGCTGGCGAAAGGTACCATTTCCCGGATGTGACGCTGGTTGCATTACTGGACGTGGACGGCGCGCTGTTTTCTGCCGATTTTCGCTCGGCAGAGC
GTTTCGCTCAGCTTTACACCCAGGTCGCCGGTCGTGCCGGGCGTGCGGGTAAACAGGGCGAAGTGCTGCAACGCACCATCCGGAACATCCTCTGTTG
CAAACGTTGCTCTATAAAGGCTACGACGCCCTTTGCCGAACAGGCGCTGGCTGAGCGGCGAATGATGCAGCTACCGCCGTGGACCAGCCATGTGATTGTGCG
TGCGGAAGATCATAACAATCAGCACGCGCCATTGTTTCCTGCAACAACATGCGTAATCTGATCCTCTCCAGCCCCACTGGCAGACGAGAACTGTGGGTTCTCG
GTCCGGTTCCGGCTCTGGCACCTAAACGTGGCGGTCGCTGGCGCTGGCAGATATTGTTGAGCACCCCTTCCCGCGTGGCGCTTGCAACACATCATTAACGGT
ACGCTGGCGCTCATCAATAACAATACCGGATTCCCGTAAGGTGAAATGGGTGCTGGATGTTGATCCGATTGAGGGTTAA
```

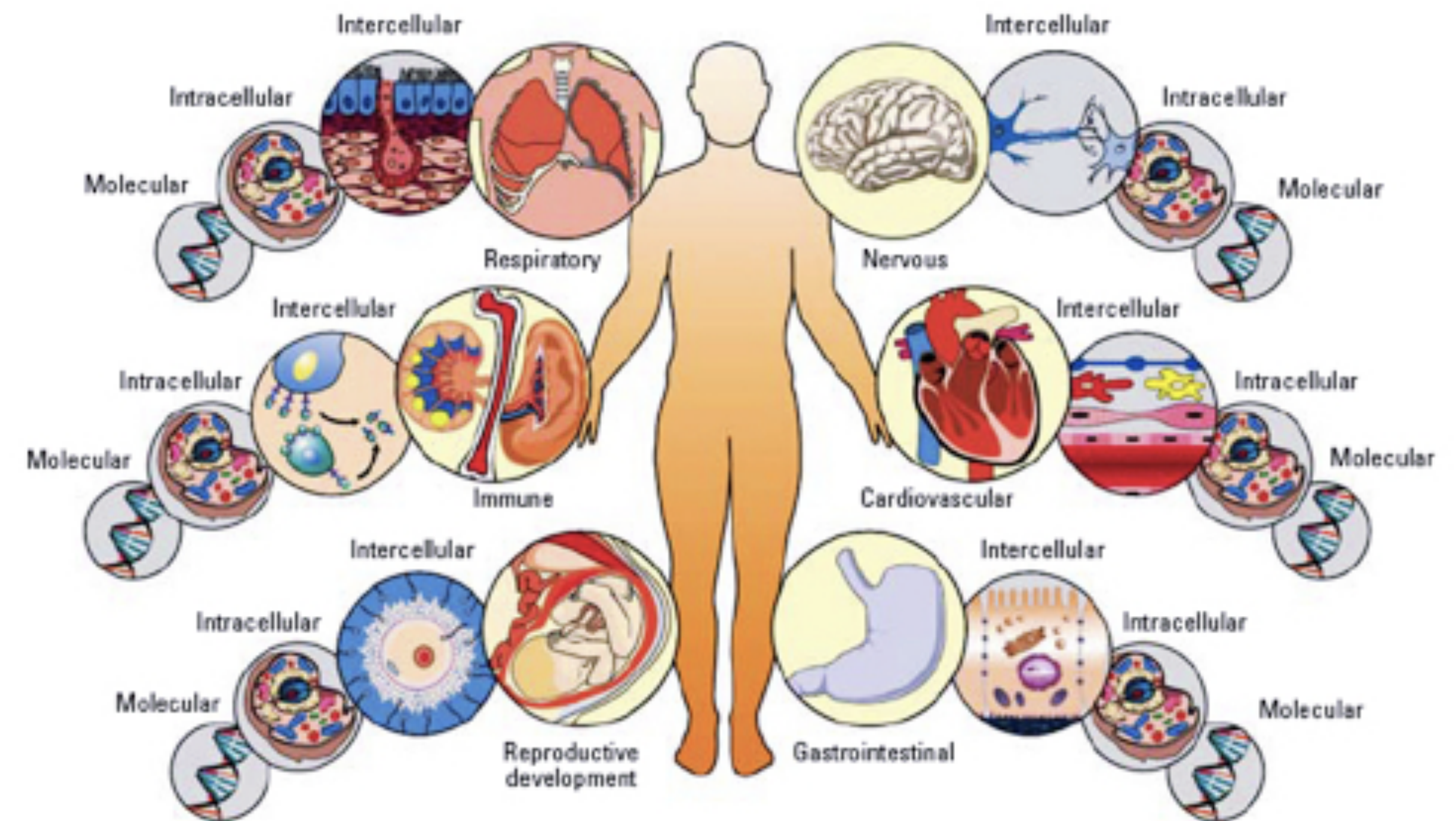


```
MPVAHVALPVPLPRTFDYLLPEGMTVKAGCRVRVPFGKQQERIG
IVVSVDASELPLNELKAVVEVLDSEPVFTHSVWRLLLWAADYY
HHPIGDVLFHALPILLRQGRPAANAPMWYWFATEQQQAVDLNSL
KRSPKQQQALAAALRQGIWRDQVATLEFNDAALQALRKKGLCDL
ASETPEFSDWRTNYAVSGERLRLNTEQATAVGAIHSAADTFSAW
LLAGVTGSGKTEVYLSVLENVLAQGKQALVMVPEIGLTPQTAR
FRERFNAPVEVLHSGLNDSERLSAWLKAKNGEAAIVIGTRSALF
TPFKNLGVIVIDEHDSSYKQQEGWRYHARDLAVYRAHSEQIPI
ILGSATPALETLCNVQQKKYRLLRLTRRAGNARPAIQHVLDLKG
QKVQAGLAPALITMRQHLQADNQVILFLNRRGFAPALLCHDCG
WIAECPRCDHYTTLHQAHHLRCHHCDSQRPVPRQCPSCGSTHL
VPVGLGTEQLEQTLAPLFPGVPISRIDRDTTSRKGALEQQLAEV
HRGGARILIGTQMLAKGHHPDVTLLVALLDVDGALFSADFRSAE
RFAQLYTQVAGRAGRAGKQGEVVLQTHHPEHPLLQTLLEYKGYDA
FAEQALAERRMMQLPPWTSHVIVRAEDHNNQHAPLFLQQLRNLI
LSSPLADEKLWVLGPVPALAPKRGGRRWRWQILLQHPSRVRLQHI
INGTLALINTIPDSRKVKWVLDVDPIEG
```



*DNA changes at very slow rates for Myrs => EVOLUTION*

```
>gi|82583714|emb|CT009680.1| Human chromosome X complete sequence
CTAACCCTAACCCTAACCCTAACCCTAACCCTCTGAAAGTGGACCTATCAGCAGGATGTGGGTG
GGAGCAGATTAGAGAATAAAAGCAGACTGCCTGAGCCAGCAGTGGCAACCCAATGGGGTCCCTTTCCATA
CTGTGGAAGCTTCGTTCTTTCACTCTTTGCAATAAATCTTGCTATTGCTCACTCTTTGGGTCCACACTGC
CTTTATGAGCTGTGACACTCACCGCAAAGGTCTGCAGCTTCACTCCTGAGCCAGTGAGACCACAACCCCA
CCAGAAAGAAGAACTCAGAACACATCTGAACATCAGAAGAAACAACTCCGGACGCGCCACCTTTAAGA
ACTGTAACACTCACCGCGAGGTTCCGCGTCTTCATTCTTGAAGTCAGTGAGACCAAGAACCCACCAATTC
CAGACACACTAGGACCCTGAGACAACCCCTAGAAGAGCACCTGGTTGATAACCCAGTTCCCATCTGGGAT
TTAGGGGACCTGGACAGCCCGGAAAATGAGCTCCTCATCTCTAACCAGTTCCCCTGTGGGGATTTAGGG
GACCAGGGACAGCCCGTTGCATGAGCCCTGGACTCTAACCAGTTCCCTTCTGGAATTTAGGGGCCCTG
GGACAGCCCTGTACATGAGCTCCTGGTCTGTAACACAGTTCCCCTGTGGGGATTTAGGGACTTGGGCCTT
CTGTCTTTGGGATCTACTCTCTATGGGCCACACAGATATGTCTTCCAACCTCCCTACACAGGGGGGACTT
CAAAGAGTGCCTTGAGCTGATCTGGTGATTGCTTTTTTGTACTGTTATTTATCTTATTCTTTTTCATTGTG
AGGTACTGATGCAACACTTTGTACGAAAAGGTCTTTCTCATCTCGGGAGTCCCCGTCTATTTGTCCCGG
TCCCTGTTAACCAGTCCCCGACAGGAGCCCTTCTGCACCTTGAGCTCTCACCCTCACCGTCCATCCA
GCCCCAGCTCTGCCTGCAACCCACCCATCCCTGGGACTCGGGCCTCCCTCTCTAGTGGTCTGGTCATCA
GGCCAGGGGCACGTGGAAGAAGCTATCGTGGAAGGGAGCAGTCATATCCCCAAAATCTGTGGTTGGTT
TACCACCACCATGGAAACCCAGGGTGGGACTCTAGTTTTCAGGTTGGAGCTGAGCCCTGTGCGGAATGAG
CTTTCCCCAGCTATGGCTTCTTGGGGCCCCCTGTGCCCTGAGCTGTGTCTCCAGCATCGGGTCCCCACCA
TGCATATGGCCCACTCAGGCACAGTGCCGCGATGGCTGCATGCGTGAGGGGGGCGCTGGGCCAGGGCTGG
GAGTCCTTTGTGTCTCATGGCCATGATTGTCCTTCCGAGTATGATATGGTGGCCAATTTCTTTTATTCTG
TCGTTTCAAGTGAAGTAAATGATGTAGAGTTCATGCAGAAAAAATACAACAAAAACCAAGGGAACATAGA
ATTGAAAACGCGTCACAGCAATGAGTTAAATAGGTAACAAATTTTCATCATTGGAAGAAAGACTTAGAGT
GCCAAAAGTGCCTCTTAAGTCTCCTTTAAAAAGTAGCAAAATTCATCCCTGAAGAAGCATCTTGGCCTTT
TTCATGTACTCAGAGTGCTGGTGAAGAACAAAGATTGCTGAAACATTATGTACCTAACAGCGTTACAGGG
TGTAATAACACACTGGAAAACCTGGTCGTTACAGTGGACATATCCAGGAAGTCCTTGCCTGAGGTTTT
CCAAGTTATGGAATTGCTTGAGATTGGAAGAGGCGATGGAGGGTACAACCTGTAATGCCCAACCTCATT
TGCTAACCCCTGTTTTTAGACTCTCCCTTTTCTTCAATCACCTAGCCTTCTTTCCACCTGAAAGGACTCTC
CCTTAAGTGAAGAACCGGACAGACTCCATCTTGGCTCTTTCACTGGCAGCCCTTCCCTCAAAGACTTAA
CTCGTGCAAGCTGACTCCCAGGACATCCGAGAATGCAATTAAGTACCAACCTACTGTGGCGAGCTACATC
CGCAGTCCCCAGGAATTCGTCCGATTGATAACGCCCAATTACCCGCGTCTATCACCTTGTAATAGTCTTA
AAGCACCTGCACCTGGAAGTGTCTTCTGTAACCATTTATCCTTTTAAACATTTTGCCTGATTTACT
TATGTAAAATTTCTTTTAACTAGACCGCCACTCCCTTTCTAAACAAAAGTATAAAAGAAAATCTAGCCCC
TTCTTTGGGACTGAGACAATTTTGAGGTTAACGCAGGGTGCCTGTAATCCTAAGGGAGGAGACCGCCACT
TCTGCTGCCCTTCCCTTCCCCACACCCCTTCTCTAGTTTATGAAACAGGGAAAAAGGGAGAAAGCAAAA
AGATAAAAAAACAGAAGTAAGATAAATAGCTAGACGACCTTGGCAGCACCACCCGGCACTGGTGGTTAA
AATAATAATAATAATAATTAACCCCTGACCTAACTACTTGTGTTATCTGTAAATTCCAGACACTGTA
TGAGGAAGCCCTGCAAAACTTTCTGTTCTGTTATCTGATGCGTGTAGCCCCAGTCACGTTCCGATGCTT
GCTCGATCTATCACGACCCCTTCAAGTGAACCCCTTAGAGTCGTAAACCCCTTAAAGGGCCAGGAATTC
GTTTTCGGGGAGCTCGGCTCTTCAAGGCCAAGTAAACCTGCCGTATCTCACCTGAGACCAACCCCAACT
ACAAAACCTCAACCTGGAATTTTCCAGGACCAAAACCCATCTATATTCTGTAACCCGAAACCTCAAAGCCT
AACCCTAACCCTAACCCTACAGTTGAGGTCCCCCGCCCTGTGGTTCCAGCTCAAGACAACCTGCCCT
TCCGTGGGTTTGCAGGCCCTCTGGTGGGGGTGGGAGCTGGGGGCCACATACAGCTCTCTGAGCTTAAGCC
```



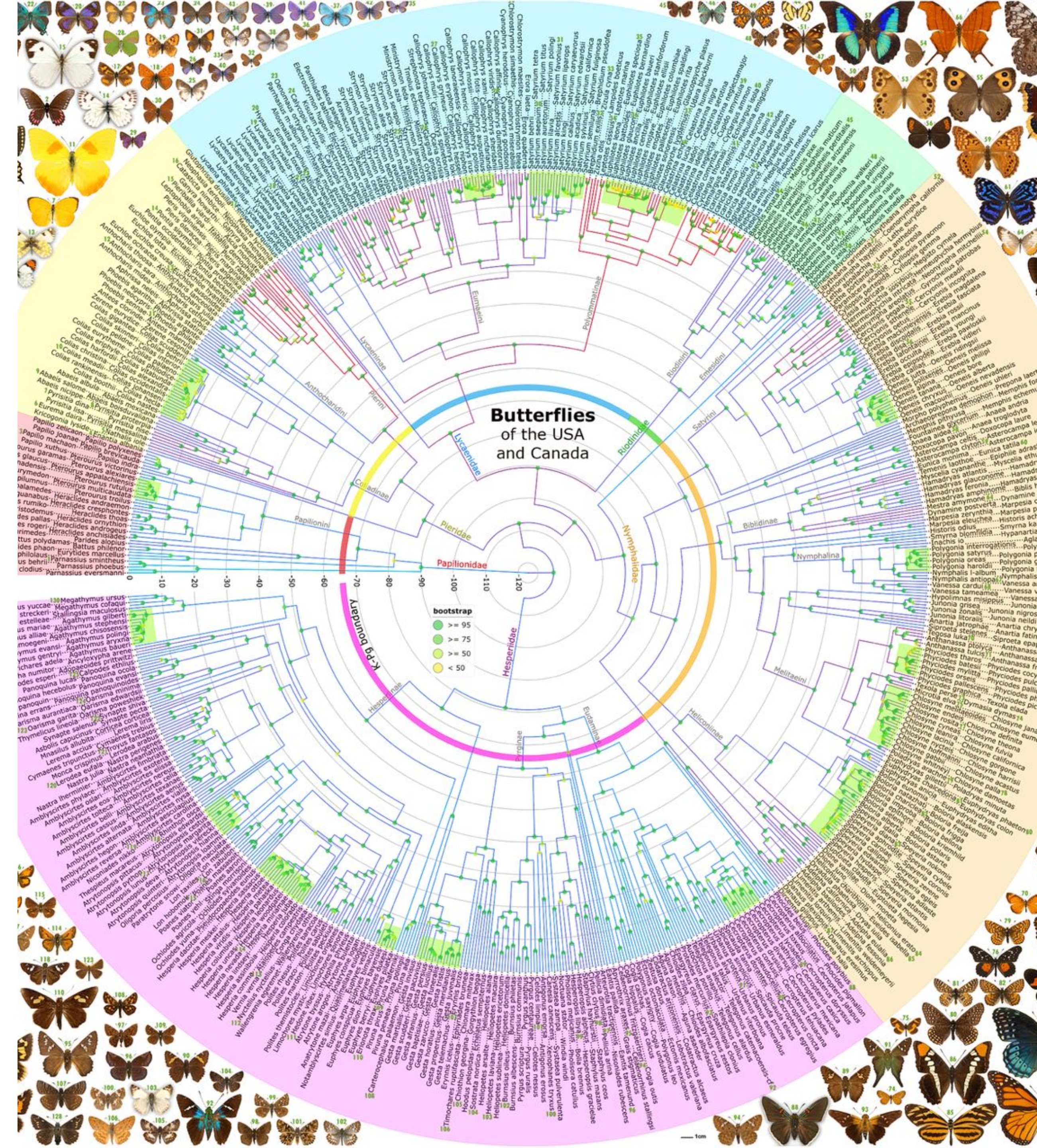
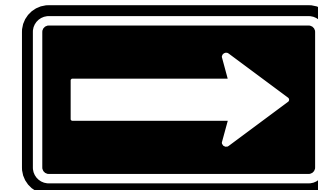
*species share many common characteristics => HOMOLOGU*



# HPC

dozens, 100s, 1000s of cores

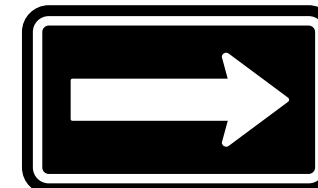
- Runs complex, parallel workloads
  - used for physics, climate, genomics
- HPC in the biosciences
  - primarily for genome science, also ecosystems
  - many comparisons, combinatorial explosion
  - time-critical (e.g. COVID-19)





*today*  
on *HYPAJA* @elixir-greece

- CPU pool VM
  - 96 CPUs
  - 1TB RAM
  - 2.4 Ghz
- Intel Xeon Cascadelake
  - >50 TB space
  - shared resource
  - project-based access etc.



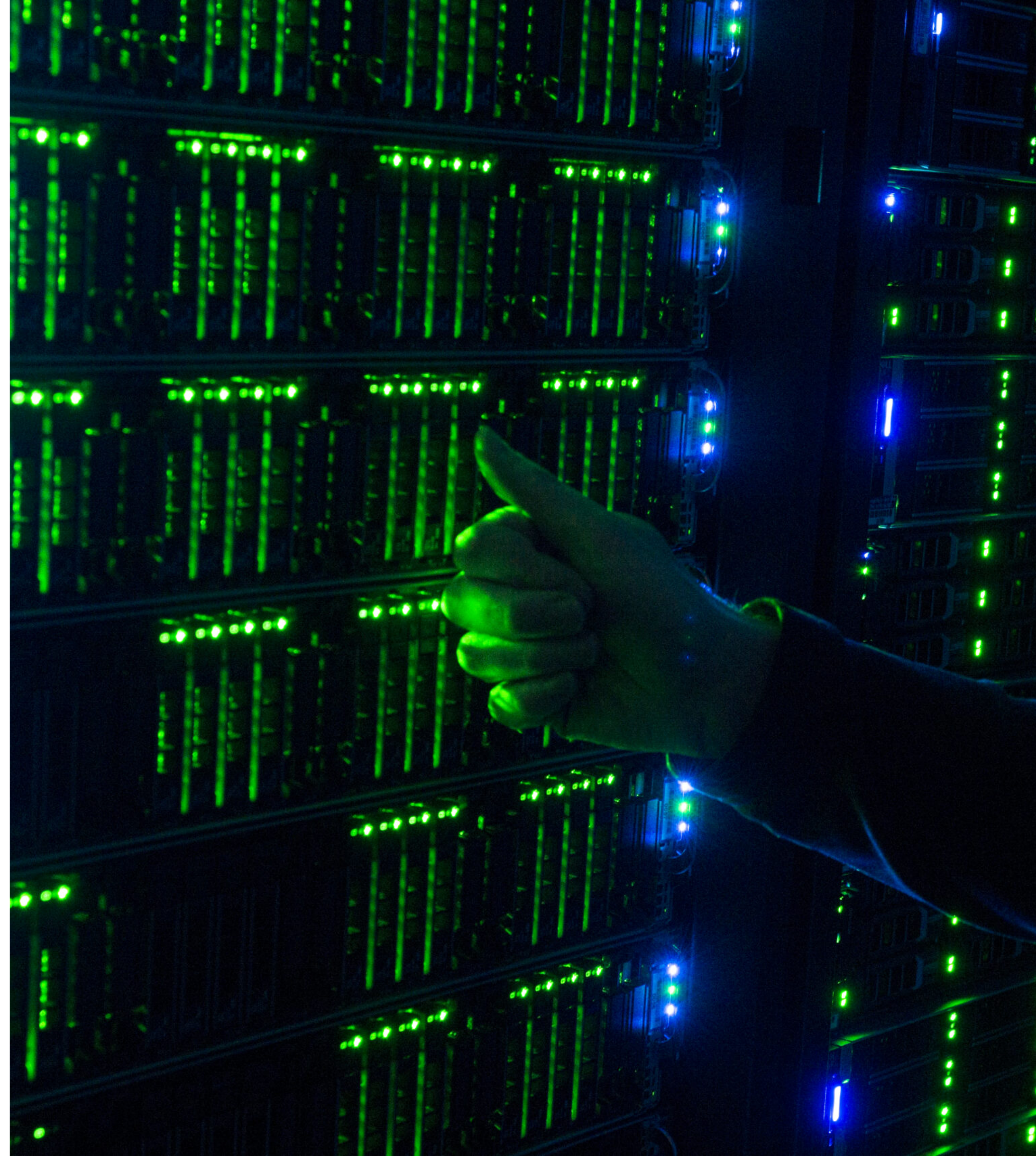
```
Architecture: x86_64
CPU op-mode(s): 32-bit, 64-bit
Byte Order: Little Endian
Address sizes: 40 bits physical, 48 bits virtual
CPU(s): 96
On-line CPU(s) list: 0-95
Thread(s) per core: 1
Core(s) per socket: 1
Socket(s): 96
NUMA node(s): 1
Vendor ID: GenuineIntel
CPU family: 6
Model: 85
Model name: Intel Xeon Processor (Cascadelake)
Stepping: 6
CPU MHz: 2400.004
BogoMIPS: 4800.00
Hypervisor vendor: KVM
Virtualization type: full
L1d cache: 3 MiB
L1i cache: 3 MiB
L2 cache: 384 MiB
L3 cache: 1.5 GiB
NUMA node0 CPU(s): 0-95
```



# *cgg toolkit*

*collection of tools, implicit*

- *Runs on HPC CPU pool*
  - *SLURM*
  - *parallel: MPI, multithreading*
- *A comparative genomics suite*
  - *input: genomes, annotations, metadata, data fusion*
  - *processing: mapping, masking, matching, clustering, visualization*
  - *output: matrices, clusters, etc.*





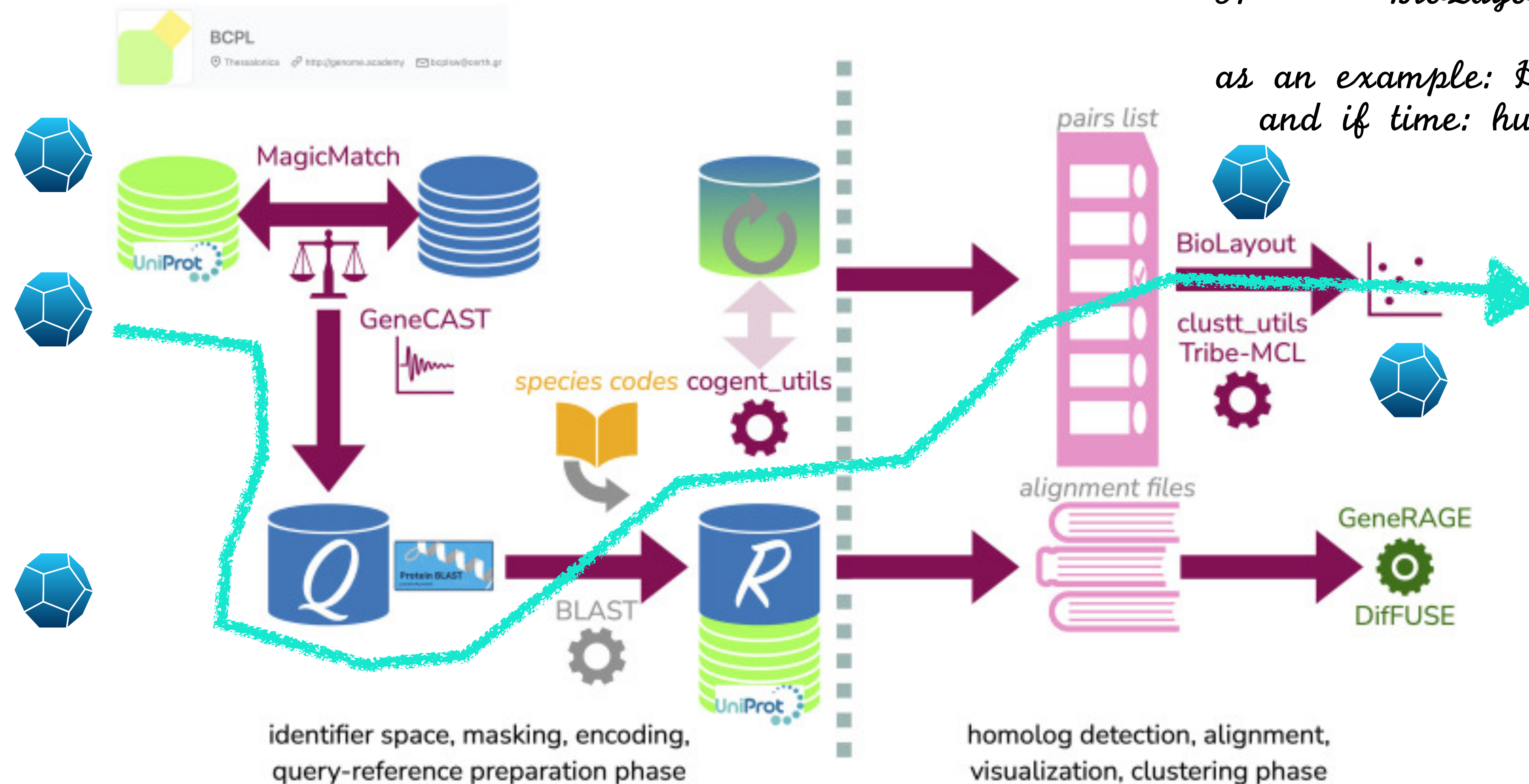


# cgg toolkit

multiple decision points for analysis

- steps today:
1. magicmatch (IDs)
  2. CAPP masking
  3. BLAST search
  4. MCL clustering
  5. BioLayout visuals

as an example: HUMAN GENOME  
and if time: human vs plant

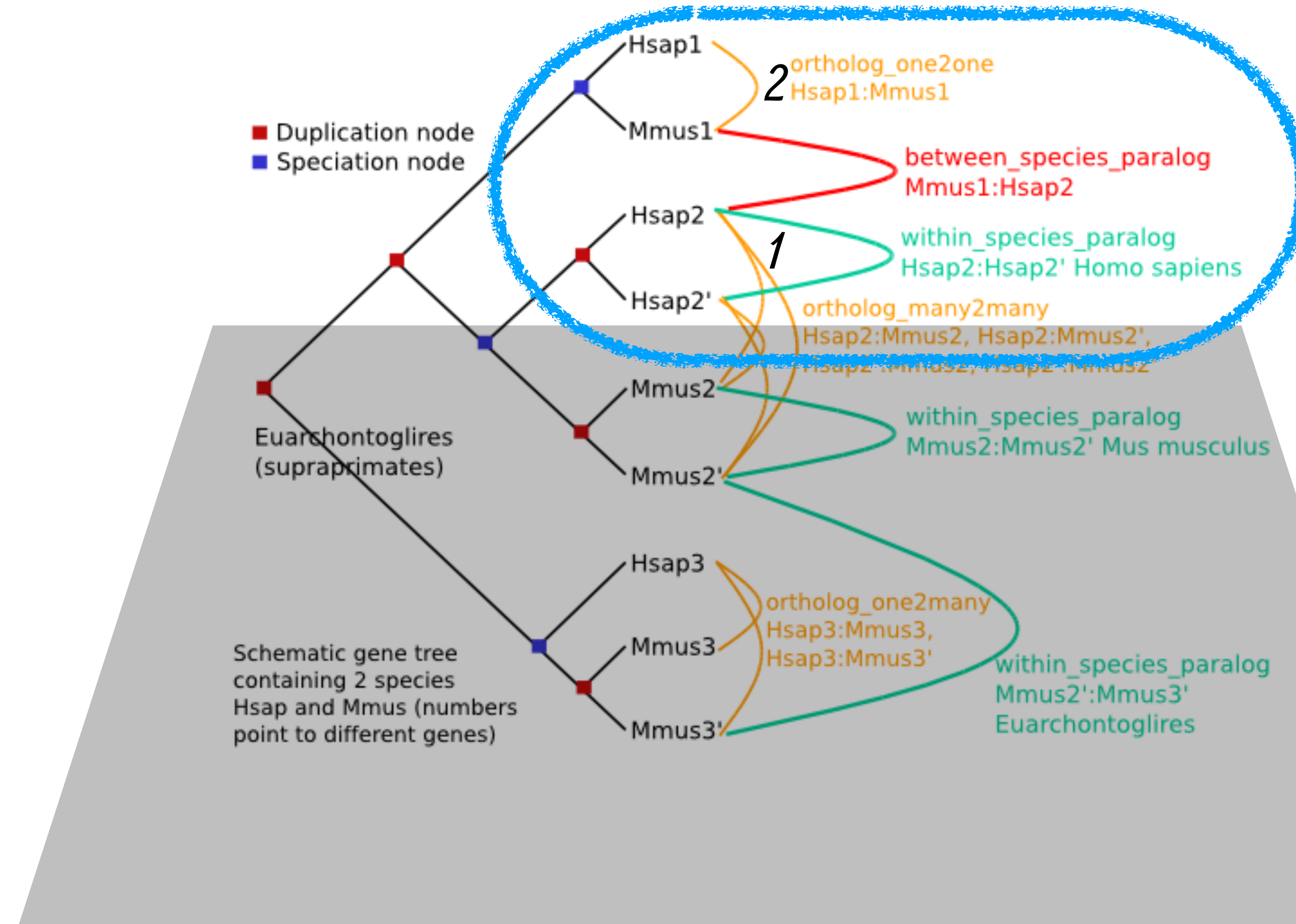




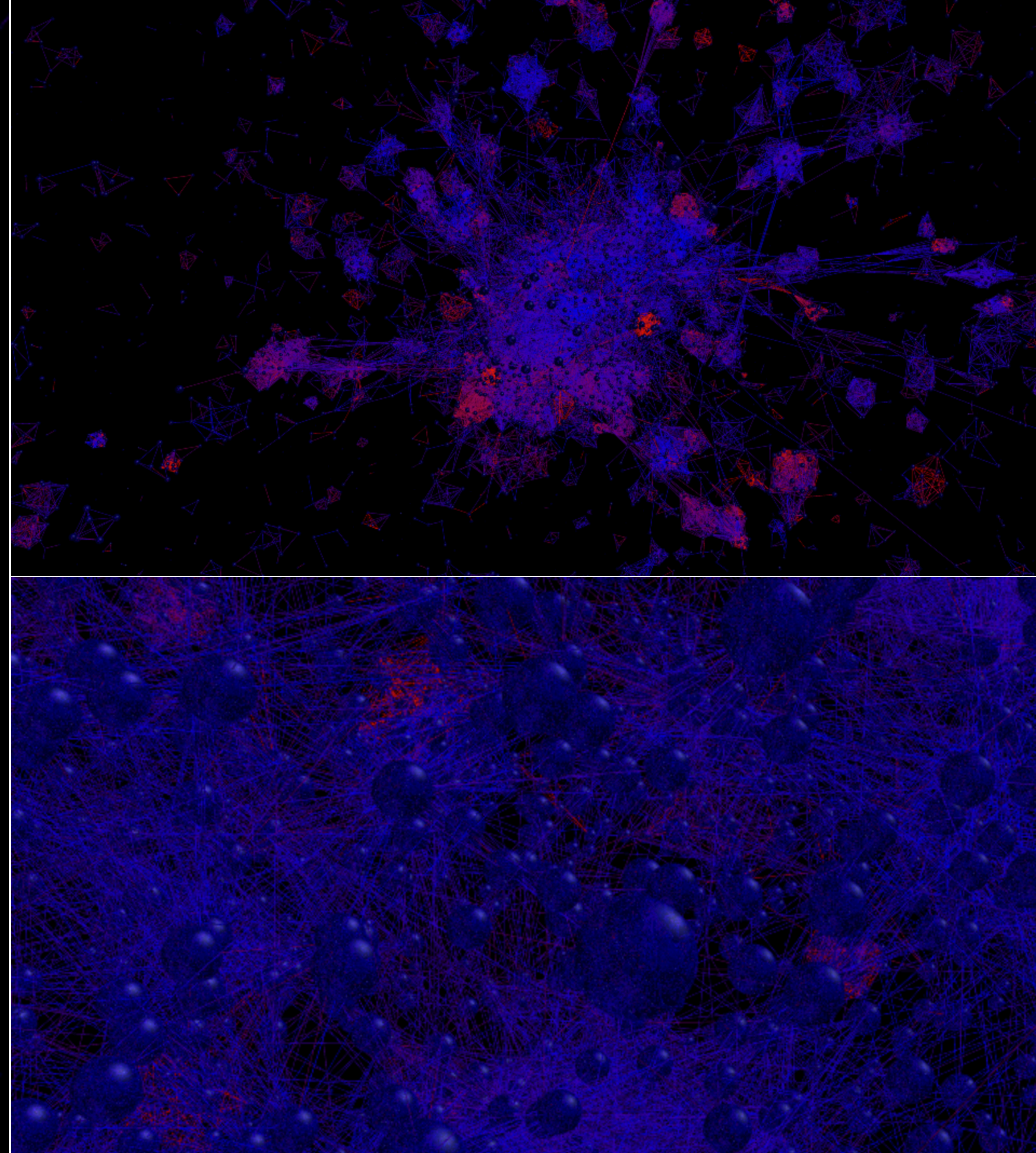
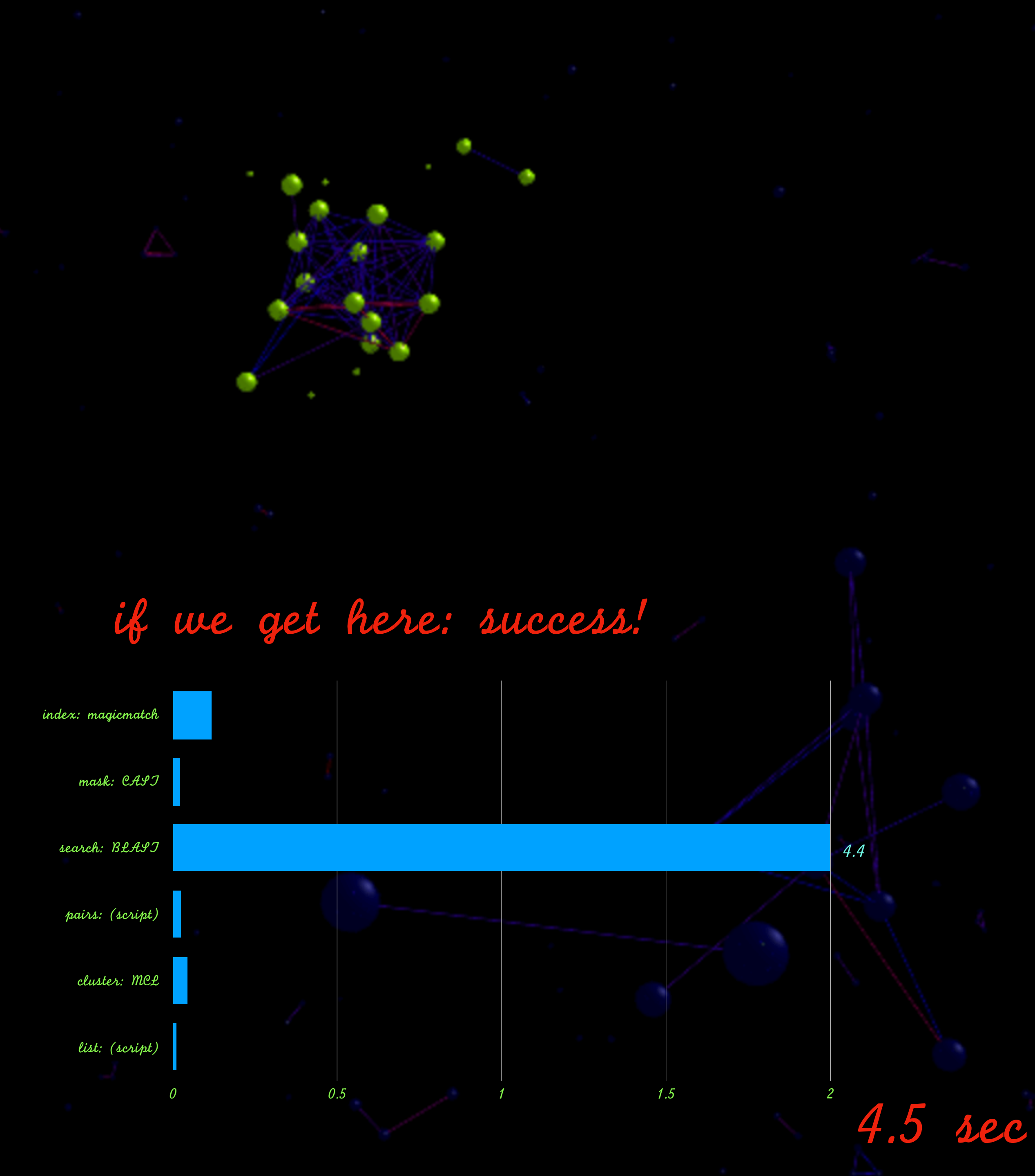
# live demo

## step-by-step

- typically automated
  - will explain config
  - will submit jobs one-by-one
- will view results
  - explain
  - discuss
  - and, if time, visualize









# conclusion(s)

and a couple of recommendations

- HPC is essential
  - for bioinformatics, much more
  - reproducibility & scalability
    - all genomes: 1bn sec = 30 yrs
- resources
  - get as much storage as you possibly can !!
    - CPU power OK, not crucial
  - prepare for **REALLY** big data !!





Q&A

<http://genome.academy/how/links>