

Bioinformatics and Evolutionary Genomics
 Evolution of Genomes, Proteomes, Networks and Complexes

Berend Snel
 Theoretical Biology and Bioinformatics
 Department of Biology
 Science Faculty
 Utrecht University

06/03/17 1

Today

- Introduction on general aims of the course and on procedural stuff
- Lecture on homology and domains (see how far we get ...)

Requests

- **very** heterogeneous with respect to previous knowledge (IBMB, GB, research projects, PhD students)
- PLEASE: interrupt / ask questions when I am going to fast, when I use jargon, when I make jumps/conclusions that to me seem obvious 100% logical, but to you are erratic; please point out my implicit assumptions regarding what everybody knows
- -> Master course ...
- Computer exercises: more experienced people help

Cancer research Agriculture Biobased economy Antibiotic resistance

HiSeq 2500 HiSeq 2500

```

    AATTGAACTTCATCGAGGCGAGCGGGGCTTTGACGAGGCTTTAGAGGGGCTTTCTTTTGG
    CTTTAGGAAGCTCAGAGGCTTGGACTCCAAAGCTTGGCGAAATTTGTGAGAGGTTTGGCGG
    GGAGTAAAGAGATGACTTCTTTGTGGCTGTGGAGAGGCTTGGGTTTGTGGGGGCTTGA
    AGTTCTTCGGGTGGAGGTAAGCTCTCAGCTCGGGGCTTCGCTTAGCTGGGCAAGTAA
    GGCGCTGGGCGGGGGGAGACTGCGCTTCTTGGCTCGGCTGACTTGTGGGAGCTT
    GGTTCATAGACTTTCGGCTCGACTCGGCTCGGGCAAGCGGCTGGGTTGCTTTGCTGG
    TTTCTGGGCTTCTCGCGGCTGGGCTTCTTGAGGCTTGGAGCTTCTTTGGAGCTT
    GGCTTTTCGGTAGGACCTCGGCTTGTATTCTGACTAAGCTTTTCTTAAAGAG
    TGACGCAATGCTAATTTTAAATCTGAGCTTCTGTGTAAAGTCTTCTCGAGTCTTGGCG
    CGCTAGAGATTTTCTTCTTGAGATTAAGCTTAAAGCAATGAGGCAATTTATGTATGG
    AAATAGTTTTAAAATCCCTCTGTTCAATCTGAGGTTTAAAGGCTTATCTGAGCTTAT
    
```

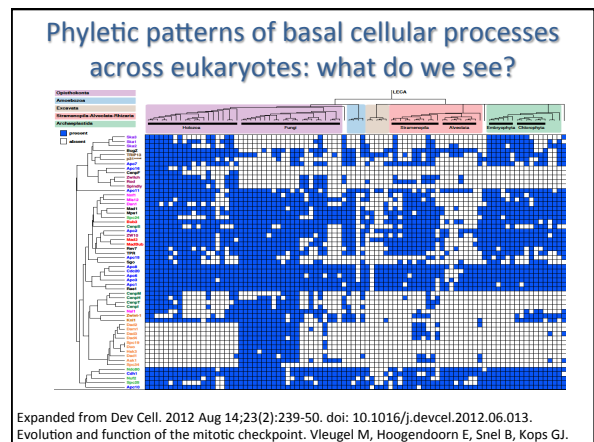
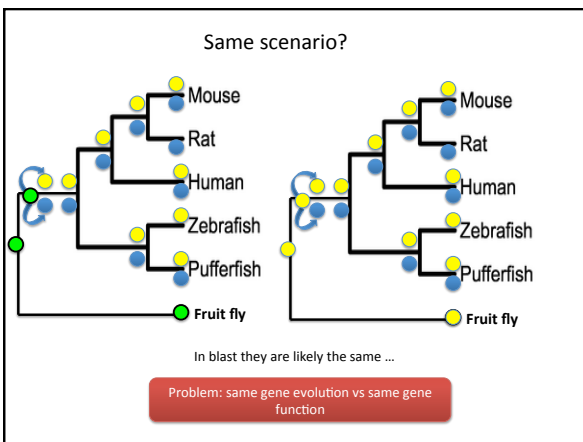
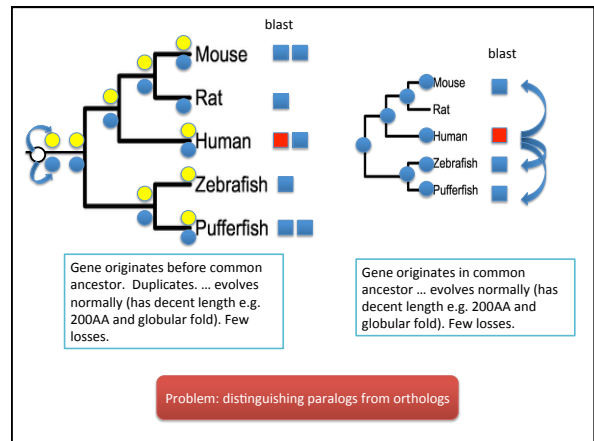
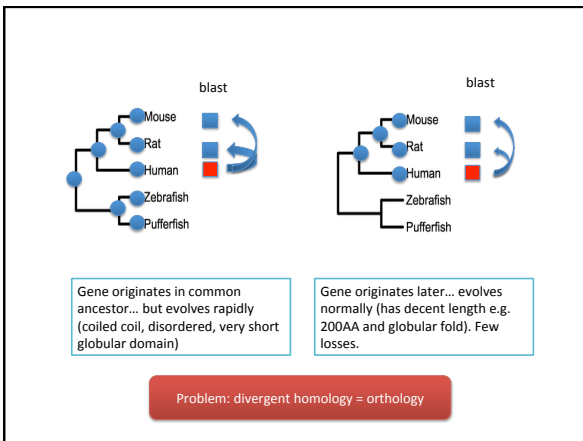
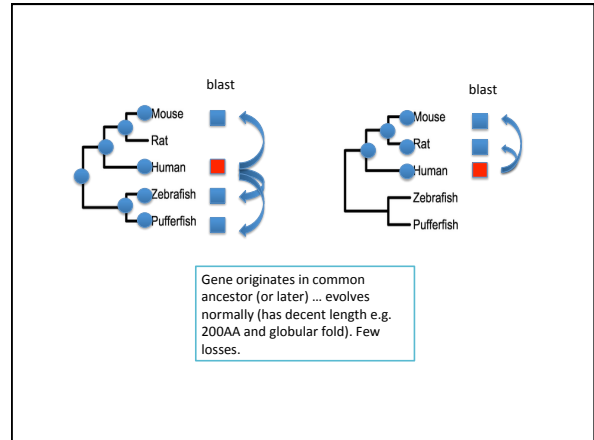
“Mapping and understanding fate of ancestral eukaryotic diversity of basal cellular processes”

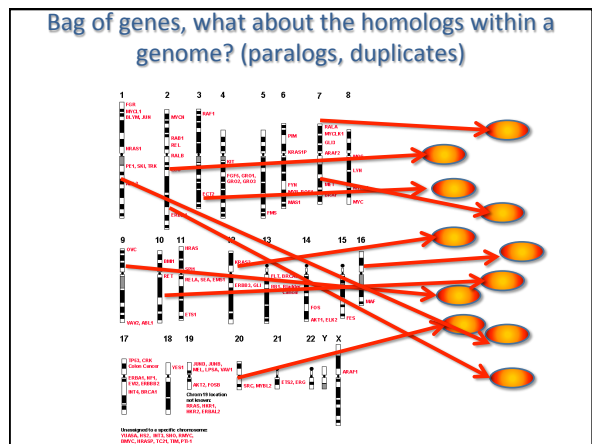
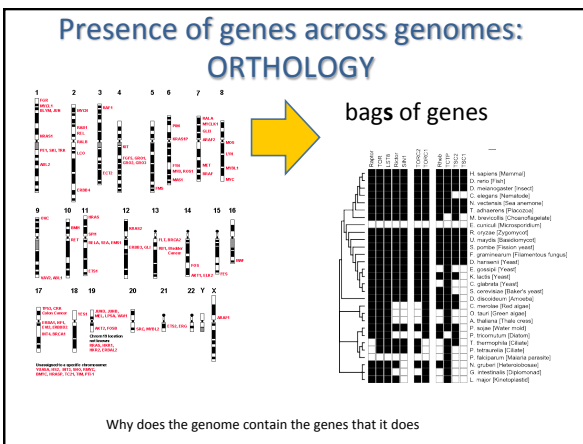
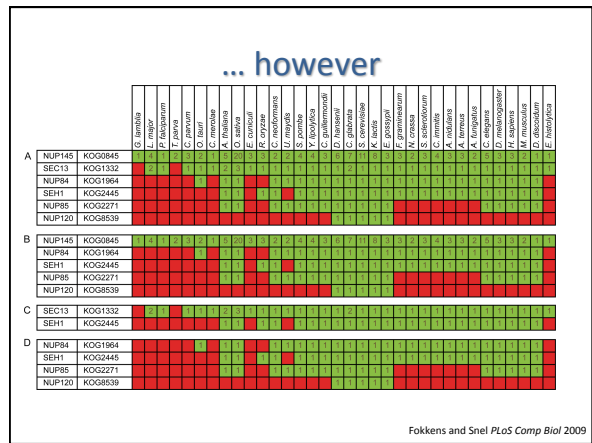
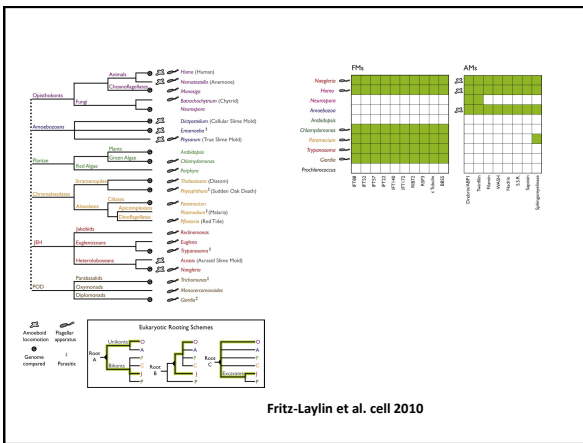
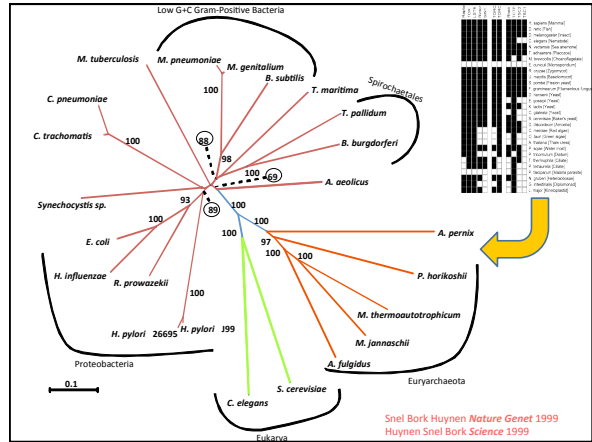
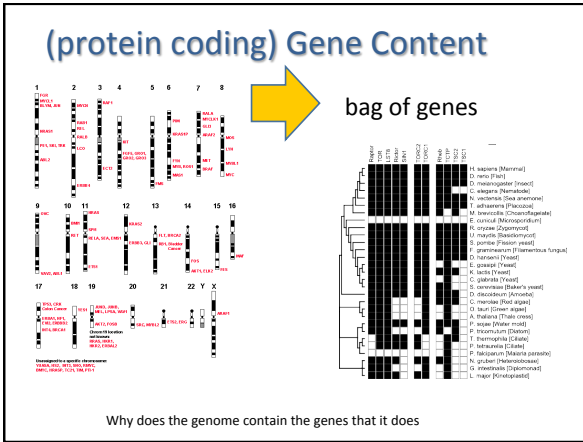
You are here

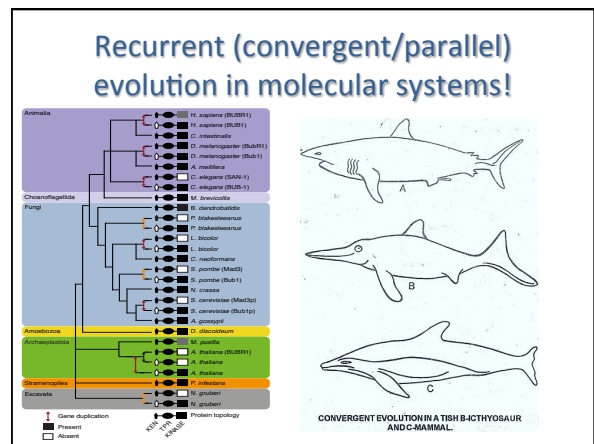
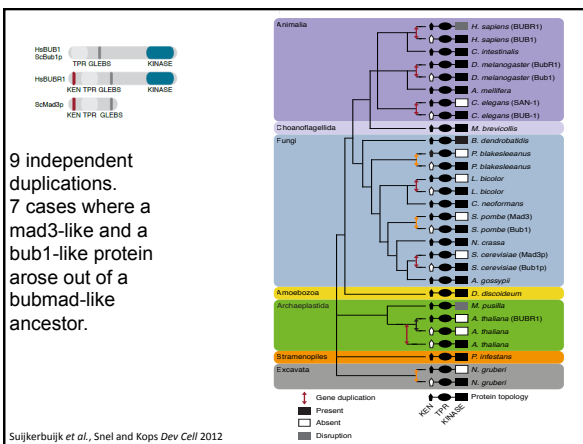
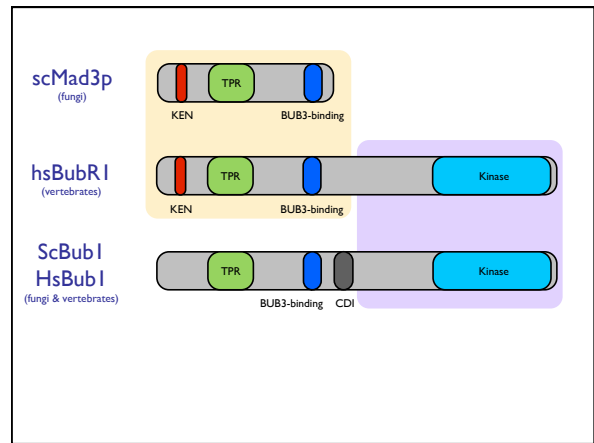
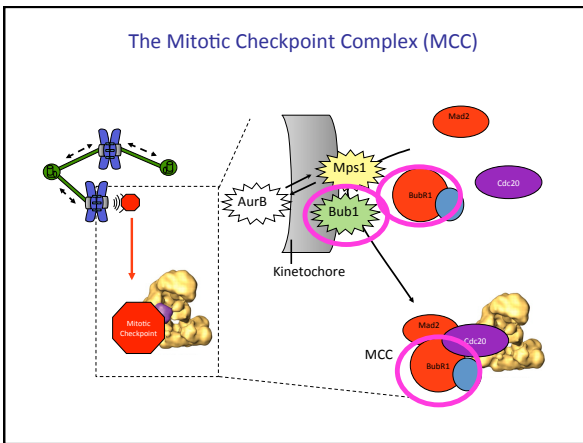
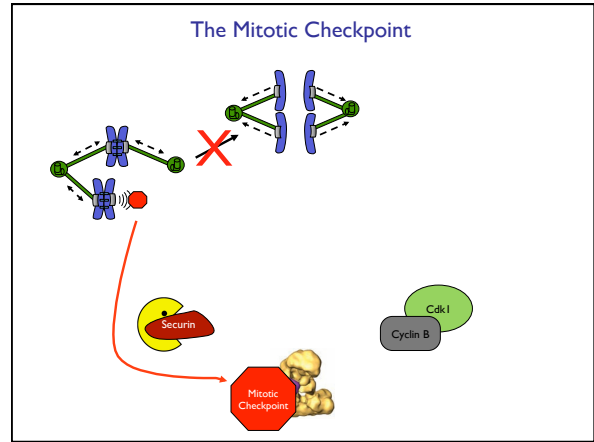
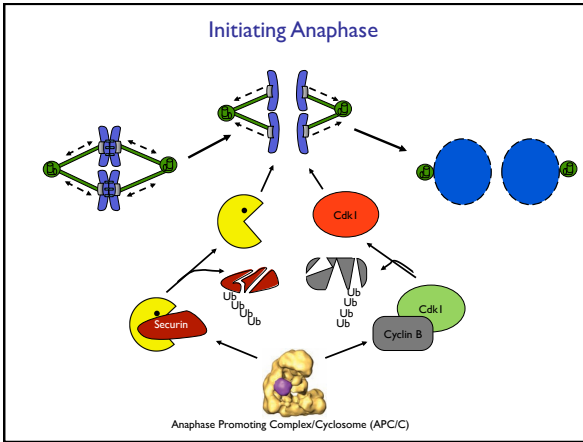
What kind of questions?

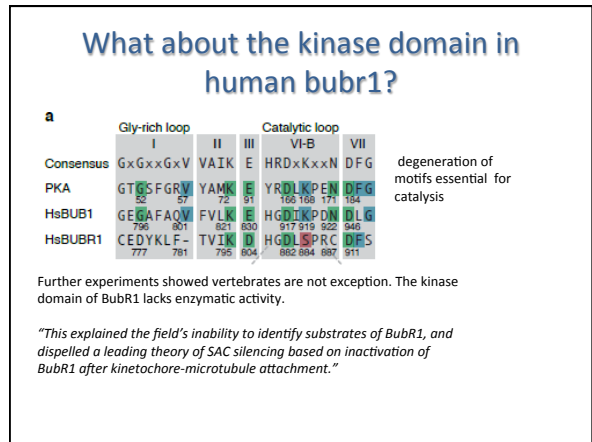
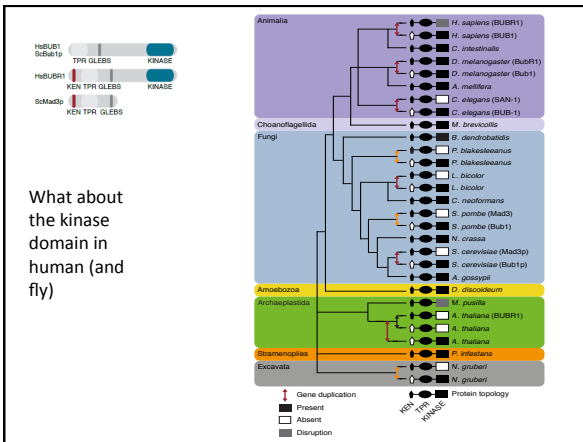
What kind of questions?

- What is the history of my gene? (what happened)
 - When was my gene invented? (what do you mean your gene? Duplication vs *de novo* gene invention)
 - How conserved is my gene? In the meaning of: in which distant species does it also occur?
 - When did this motif arise?
- “± which other genomes contain my gene”
- “What is the same thing in different organisms and how did it evolve ... “







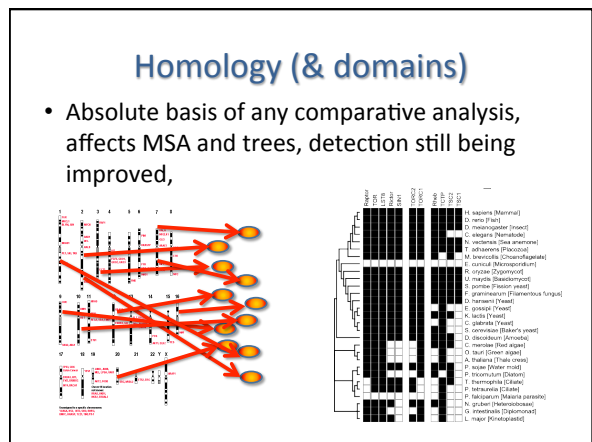


This course

- I want study evolution of genomes pathways and networks, so that is why I study gene/protein evolution
- At the end, able to analyse your own protein
- Understanding that many bioinformatic challenges are a mix of conceptual and technical problems (e.g. why orthology is such an incredibly persistent problem)
- "what you should ~know" in order to this kind of research
- Topics are interrelated
 - e.g. orthology already in homology lecture but proper explanation a day after
 - e.g. that trees can be used to time a duplication to eukaryogenesis but proper discussion of eukaryogenesis has its own lecture

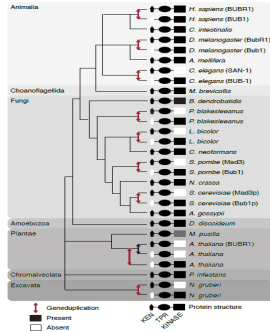
Small scale & Large Scale

- How did my protein, complex, pathway evolve? (collaborations)(COO, mini project)
- Large scale, how do genome, networks and complexes evolve (context/expectation, bioinformatics senior authorships)(paper discussions)
- What can we infer about eukaryogenesis?



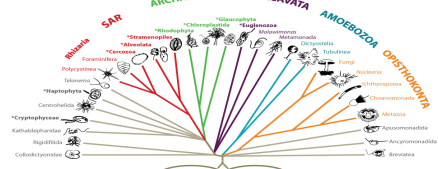
Gene Phylogeny & Orthology

- How do we get such trees and how do we interpret them
- Trees reveal some of the most important genome evolution processes: LGT, duplication, loss,



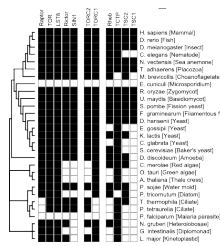
(Eukaryotic) tree of life & eukaryogenesis

- Which genome to include. What does an absence mean?
- Essential for interpreting gene trees:
 - Knowing (at least the outline) by heart >>> having to look it up
- With regards to evolutionary signaling cell biology (kinases, smallGTPases etc.)the diversity in present day genomes is staggering and dwarfs e.g. human-fruit fly difference

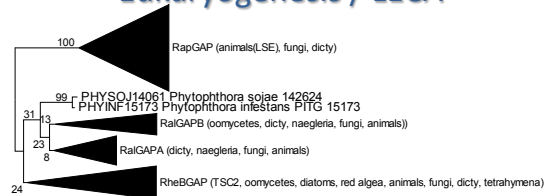


Large scale orthology

- Needed to move beyond anecdotes, but difficult to get



Eukaryogenesis / LECA



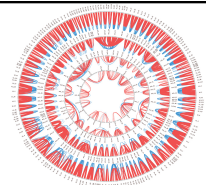
- Biological topic, eukaryogenesis / LECA for which these types of analyses are telling us a lot. But it also impacts a lot of things we do: we see it back in gene trees and it impacts getting orthologous groups across eukaryotes.

Gene content evolution



- Fundamental level of genome evolution
- Gene invention -> inability to detect homologs vs real lack of homologs does not simply mean novel gene
- Evolutionary modules?
- Trying to move large scale but remember the pitfalls

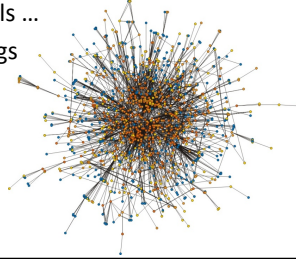
Whole Genome Duplication (WGD)



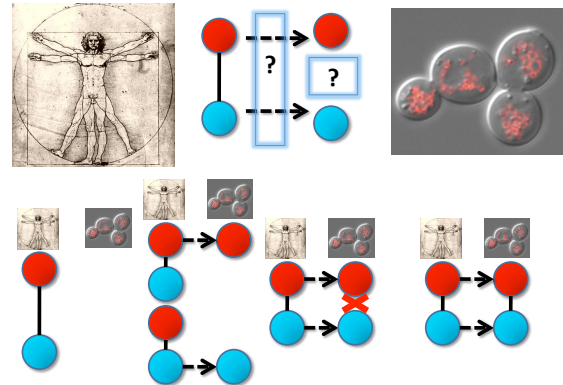
- Like LECA, WGD is important biology for which bioinf needed to research but which also impacts our data
- And which is welcome source of information for our analyses (Lidija, bubmad): independent and reliable reconstruction of the history of part of the history of genes

Using HTP data to study evolution of networks / complexes

- Is the number of conserved interactions between e.g. yeast and human 10% or 95%???
- On top of all the genome analysis pitfalls also all the HTP data pitfalls ...
- Duplicates vs orthologs



What can happen to an interaction in evolution



Techniques AND biology

- Detective/forensics vs concepts; Large scale biology vs small-scale biology; Bioinformatics biology vs data/technique problems;
- A lot like police investigation ... and less like Nobel prize winning physics ...
- Anything goes in genome evolution; many processes often entangled (i.e. google subneofunctionalization)

Practical stuff

- Schedule
- BYOD
 - Tuesday March 7
 - Tuesday March 14
- Literature discussion (have you already read Zmasek et al?)
 - You should have read the papers in depth before the discussion
 - I will shortly introduce and then invite people to discuss figures / pieces of the results
 - This + participation in the discussion is 20% of grade
- Mini project, let me first explain some bioinformatics ... than tomorrow afternoon let's discuss it & pick proteins

Computer Exercises

- Mostly use of web resources.
- Computer exercises for some topics many others more difficult (i.e. evolution of interaction networks based on HTP analysis).
- Ask help from fellow students.
- Should tie strongly into mini-projects
- (I am slightly afraid the data bases are getting unwieldy w.r.t. number of genomes ... searches very slow ... you need to already know the ToL to pick relevant species)

Requests

- **very** heterogeneous with respect to previous knowledge (IBMB, GB, research projects, PhD students)
- PLEASE: interrupt / ask questions when I am going to fast, when I use jargon, when I make jumps/conclusions that to me seem obvious 100% logical, but to your are erratic; please point out my implicit assumptions regarding what everybody knows
- Computer exercises: more experienced people help
- And also apologies for some redundancy

Mini project

- **The protein.**
- **What does my protein look like (protein topology)**
- Optional: Size of the (super)family in the genome you're sequence is from
- **Homologs across tree of life**
- **Tree of relevant sequences in diverse genomes**
- **Orthologs in genomes from your tree. (or from homology searches)**
- Optional: Does your protein or any of its orthologs in other species have Whole Genome Duplicates (WGD)/ Ohnologs?
- Optional: Point of invention of the eukaryotic orthologous group your protein belongs to.
- Optional: interactions of proteins in your tree according to biogrid
- Optional: orthology of interactors of your proteins according to biogrid and an automatic orthology database such as. E.g. panther.

Mini project

- Species tree, you really get to know the outline if you are using the ToL to describe the evolution of a protein. Similarly for e.g. smart/pfam etc.
- Students are often finished too long after the course ... for your own benefit try to prevent that
- Some students get stuck on is what they find novel. It does **not** have to be novel! Just describe what you find!

Mini project / Molecular evolution is recursive / iterative 3: generalized

- To study the evolution of a gene you need a model / framework of the evolution of the gene, but to get an idea of a proper framework / model of the evolution of a gene, you need the need to study the evolution of a gene
- Thus: heuristics & build on previous results. Start from stuff you trust (alignment of highly identical sequences), and/or only use the general but flawed overview (e.g. guide tree). Then iterate
- Not yet so automatically solved for evolutionary history of a gene and its homologs as it is for other case ...