

Evolution of Protein Post-translational Regulation

Reviews in Computational Biology

Pedro Beltrao (EMBL-EBI)

On the menu



1. Very brief introduction to post-translational modifications (PTMs)
2. Evolution of PTM regulators
3. Evolution of PTM sites and interactions
4. Suggested focus and structure for the comp bio review

Bioinformatics / Comp Bio

Biology

Problems

Systems

Questions

Interpretation

Observations

Datasets

Bioinformatics

Visualization

Statistics

Databases

Probability theory

Dynamical Systems

Processing Power

Matrices

Mathematical objects

Programming

Infrastructure

Numbers

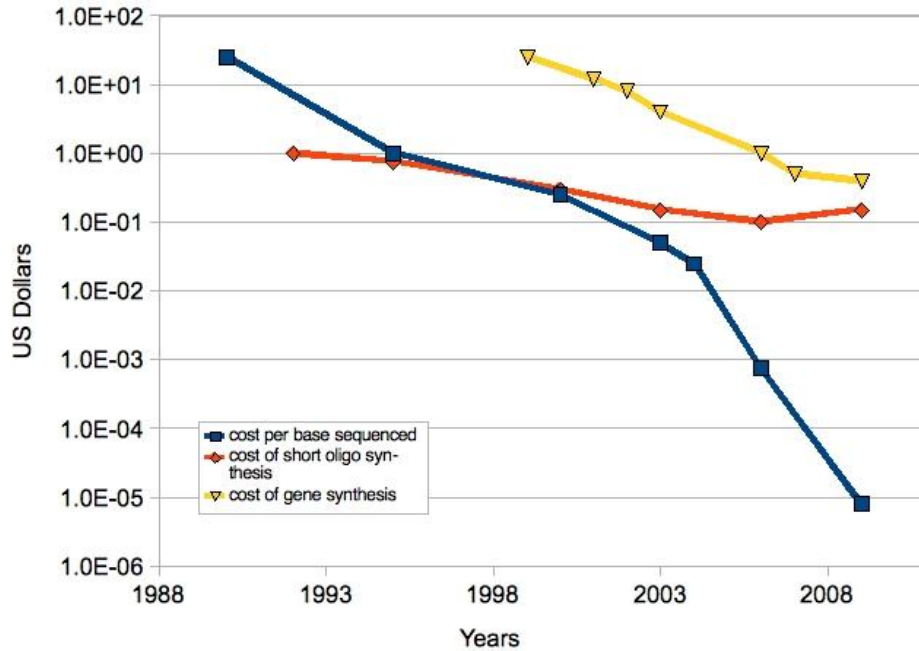
Math

Computer Science (CS)
Information Technology (IT)

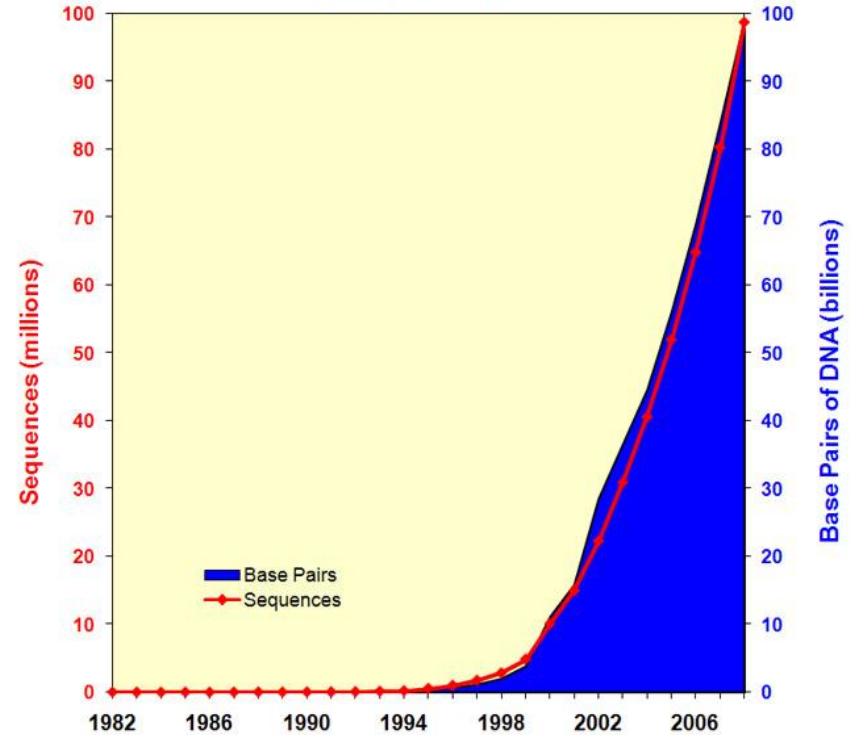
Exponential growth of data production

Cost Per Base of DNA Sequencing and Synthesis

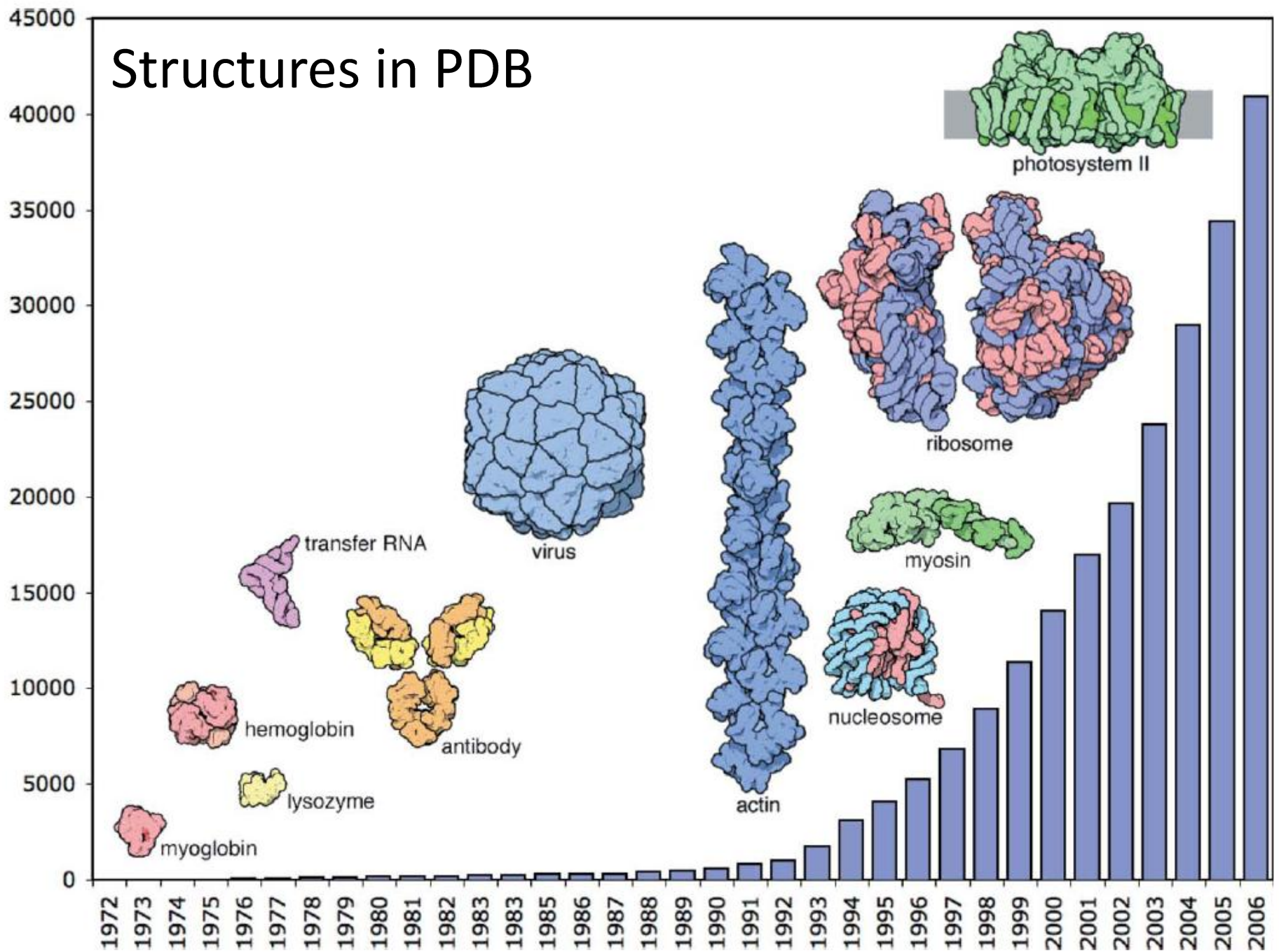
Rob Carlson, September 2009, www.synthesis.cc



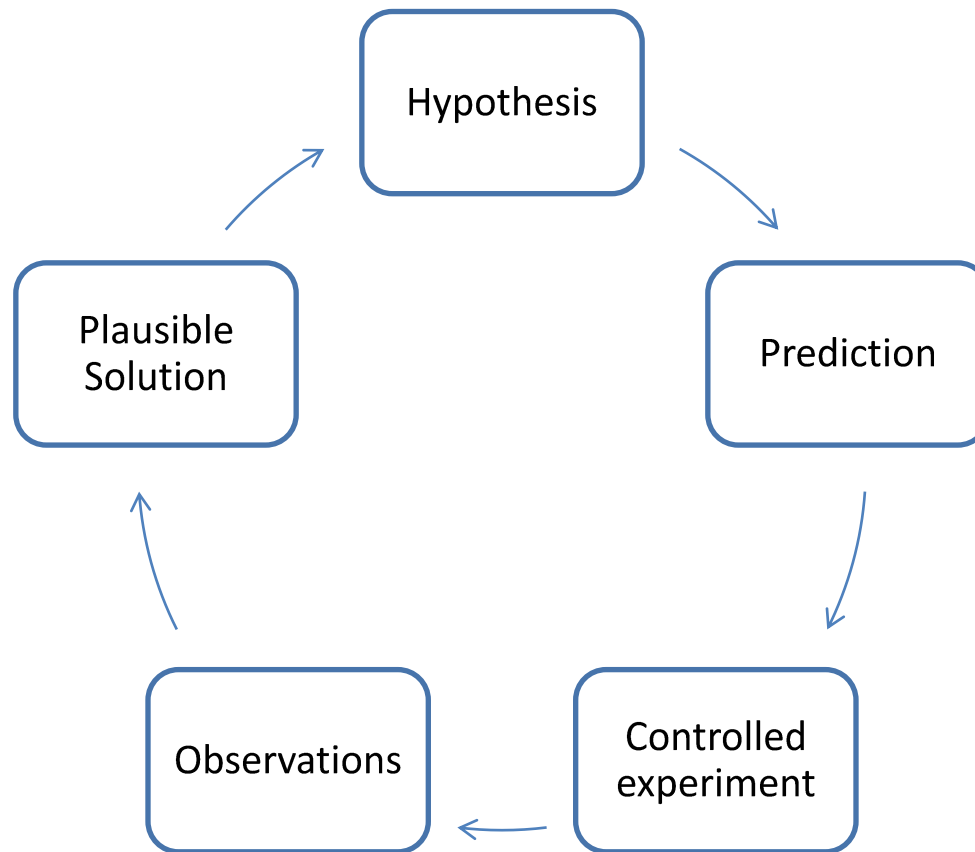
Growth of GenBank (1982 - 2008)



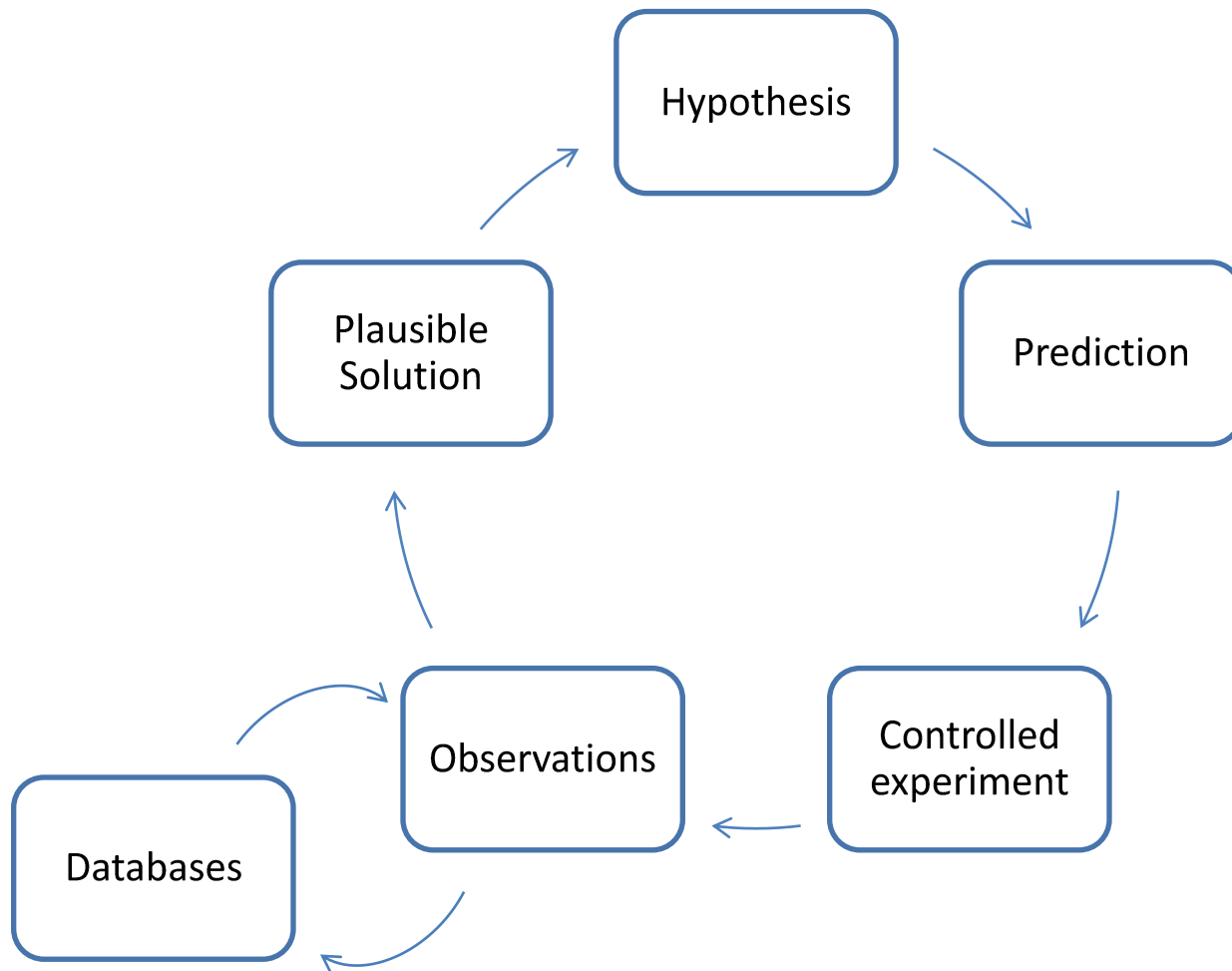
Structures in PDB



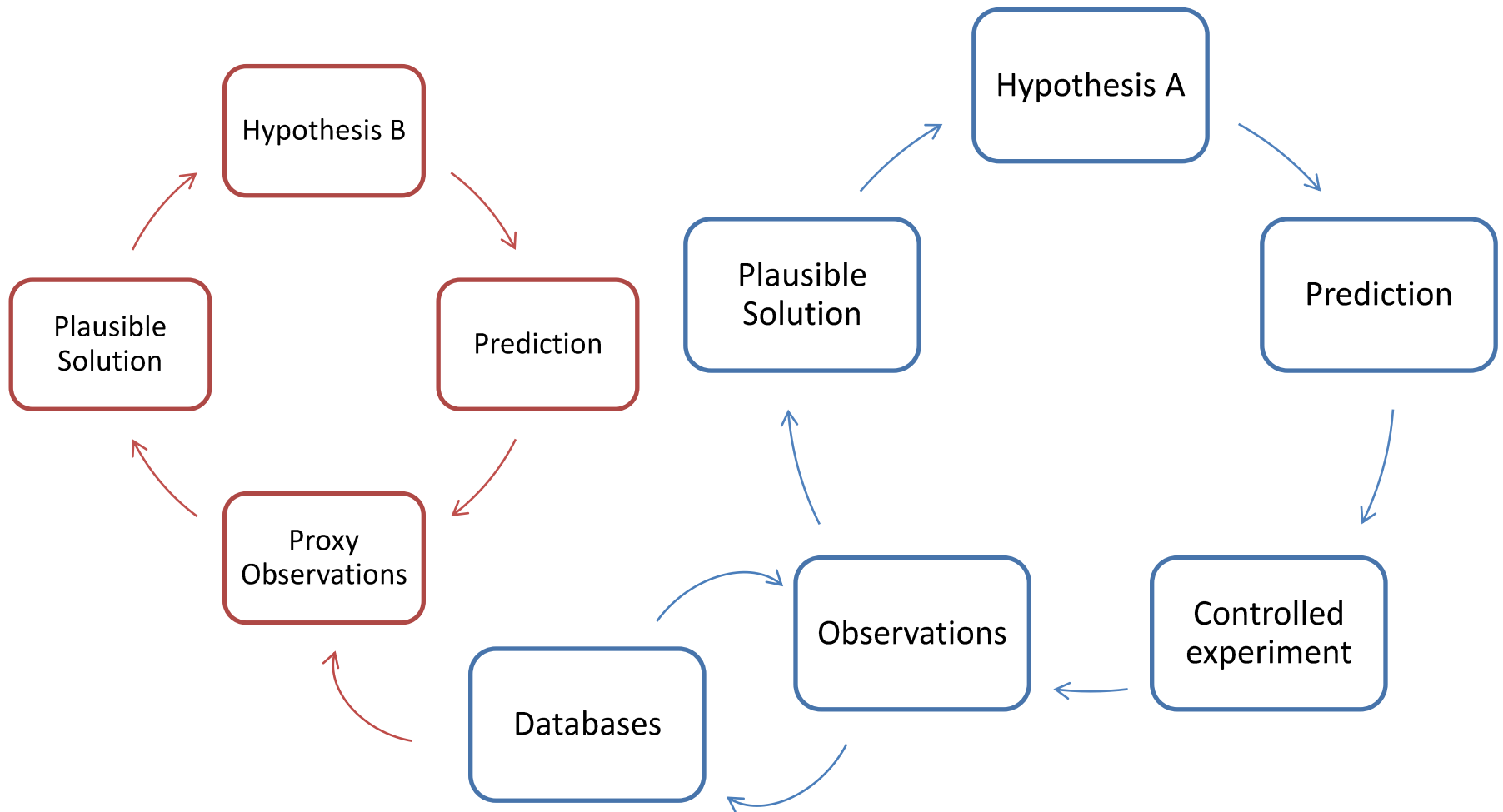
The scientific process



The scientific process



The scientific process





what am I thi

what am i thinking right now

what am i thinking

what am i thinking game

what am i thinking quiz

what am i thinking about

what am i thinking of right now

what am i thinking right now quiz

what am i thinking quotes

what am i thinking right now game

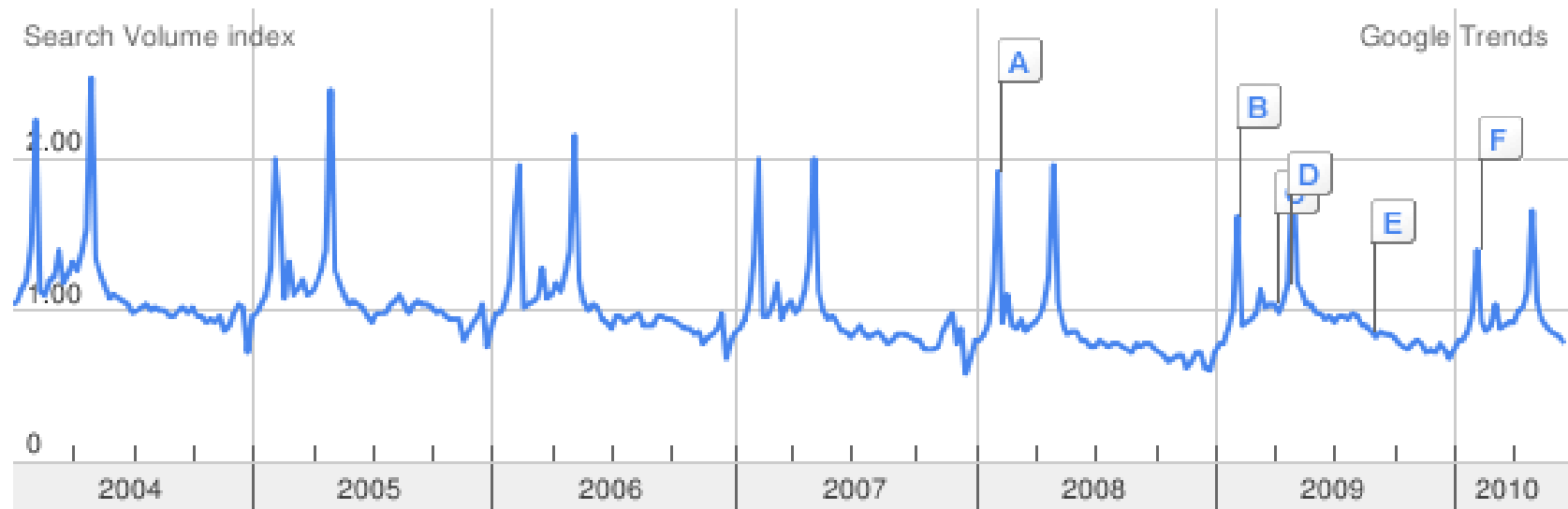
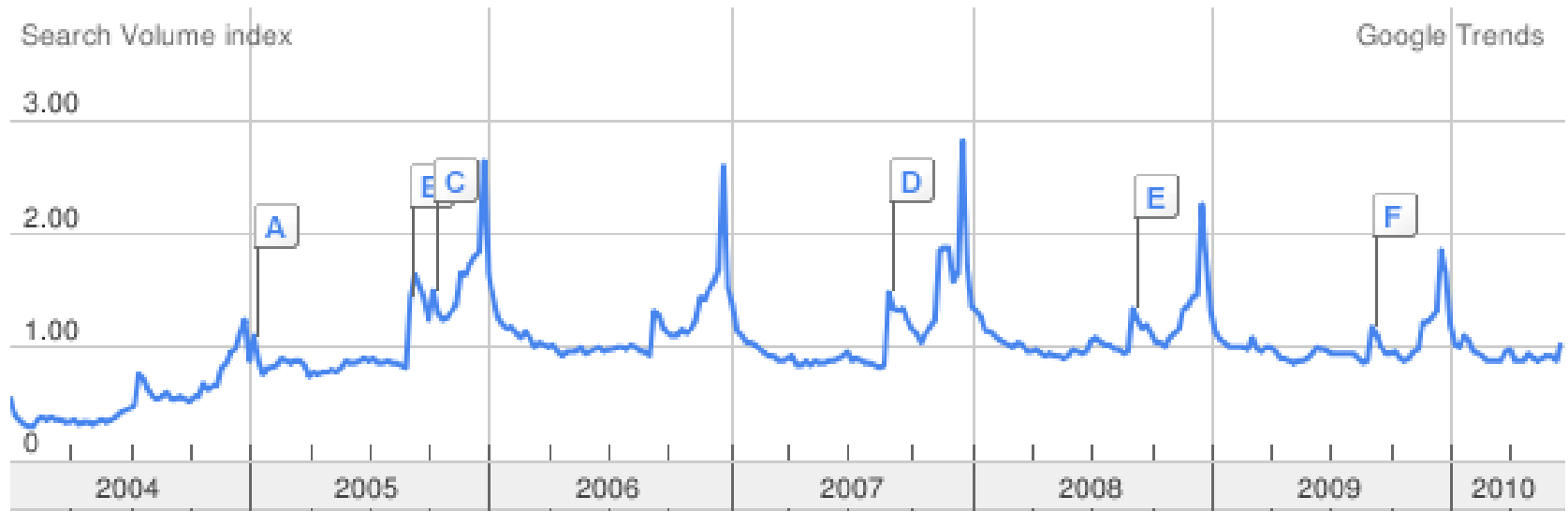
what am i thinking right now drake lyrics

Google Search I'm Feeling Lucky

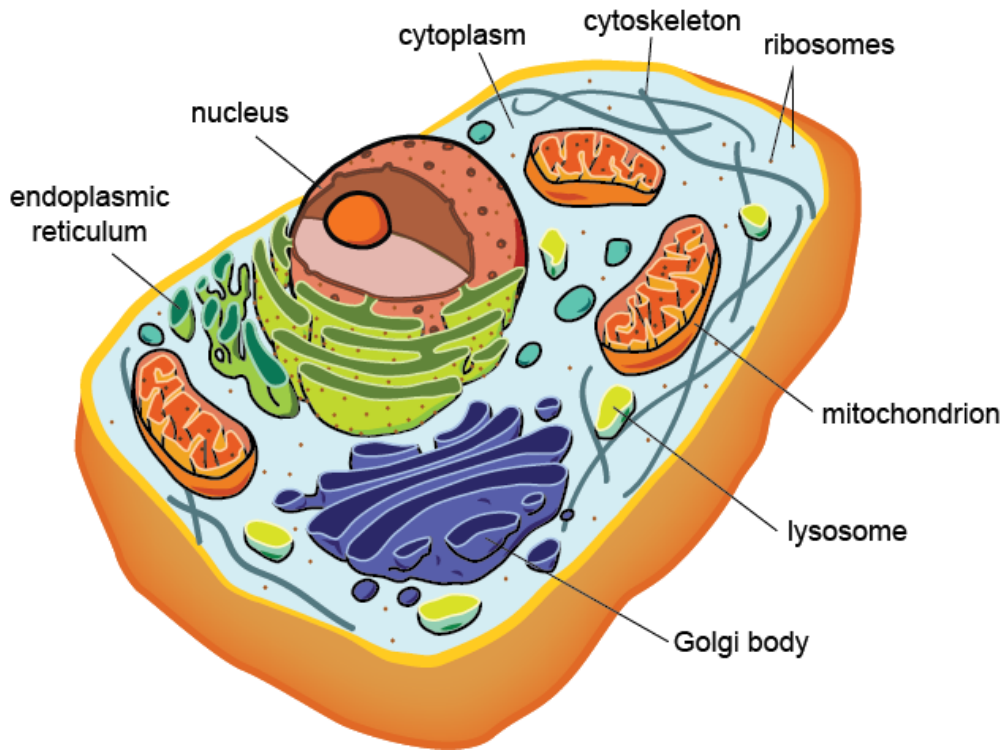
[Advanced Search](#)
[Language Tools](#)

With the added benefit of having transferable skills applicable to other “big data” problems

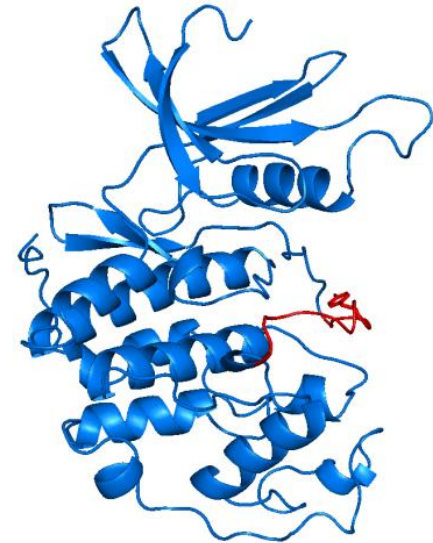
Flowers and iPods



The life of a cell



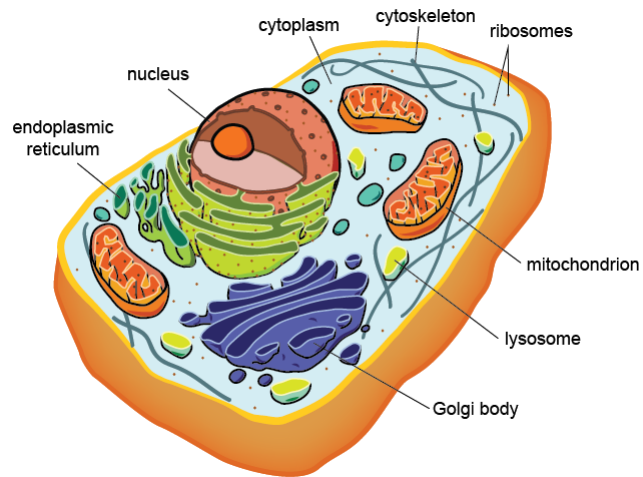
Average size of a cell is $\sim 10^{-6}$ m



Average size of a protein is $\sim 10^{-8}$ m

Some proteins are present with just a few copies inside the cell (transcription factors), some up to 10^6 copies (ribosome)

The life of a cell



What impressions did you get about how a cell works ?

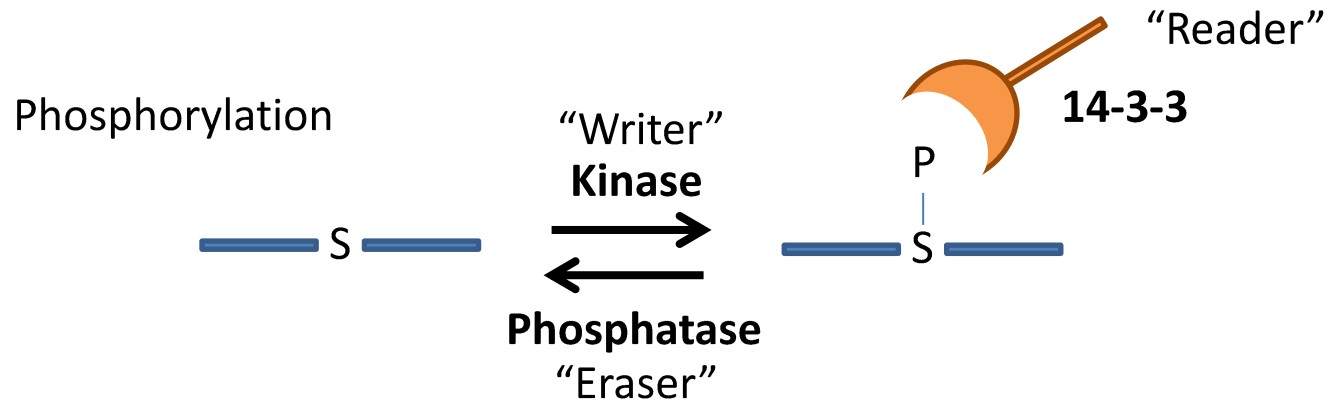
Protein Post-Translational Modifications (PTMs)

What are PTMs ?

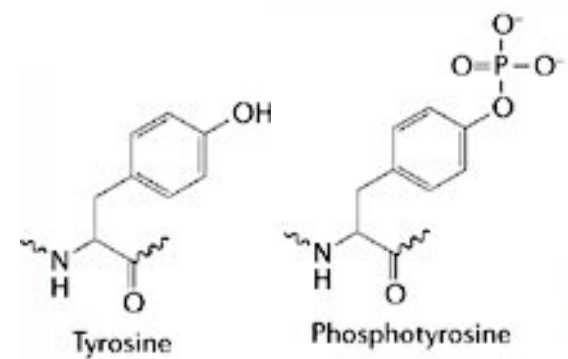
Can you give me examples of types of PTMs ?

Why would the cell invent PTMs ?

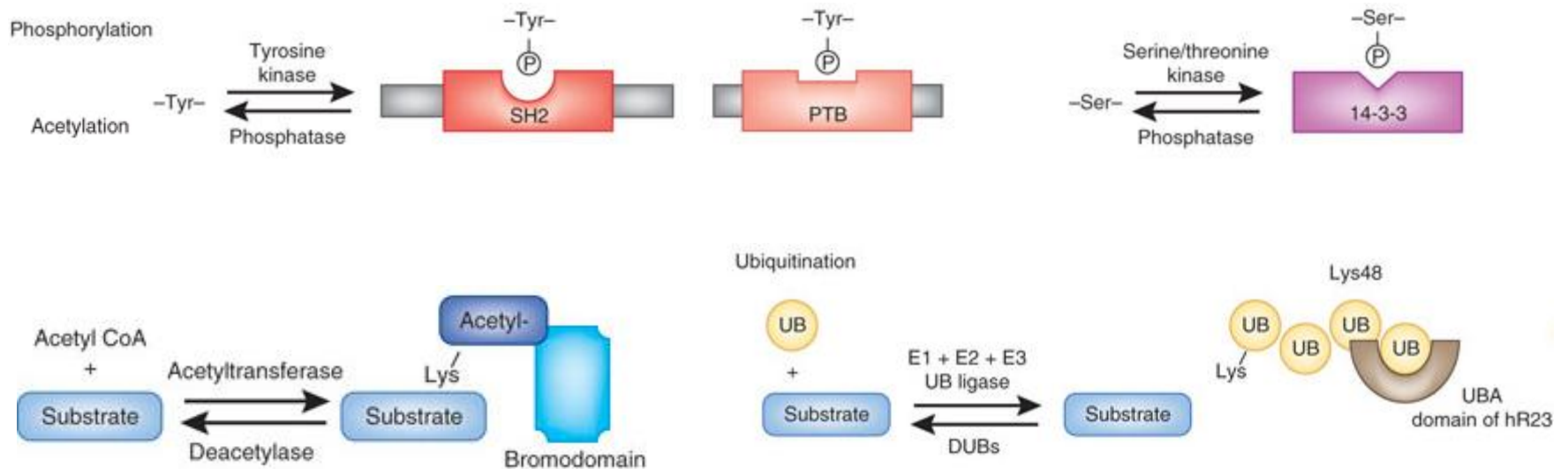
Protein Post-Translational Modifications (PTMs)



- Several amino-acids (AAs) can be phosphorylated by different types of kinases. The most commonly modified are Serine (S), Threonine (T) and Tyrosine (Y)
- The phosphorylation increases the size of the amino-acid and adds negative charge.

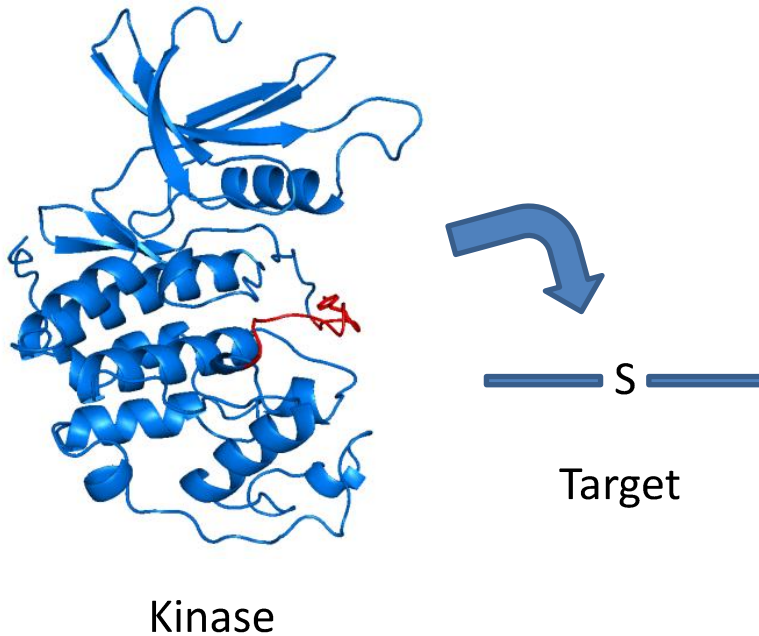


Protein Post-Translational Modifications (PTMs)



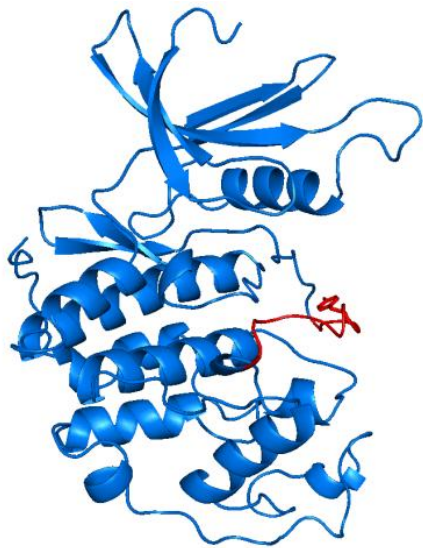
Besides regulating interactions with “reader” domains. PTMs can control function via other ways. Can you think of examples ?

Specificity of PTM regulation

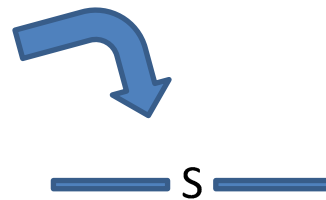


There are many Serines in all the proteins inside the cell. How does a kinase know what Serines to phosphorylate ?

Specificity of PTM regulation



Kinase



Target

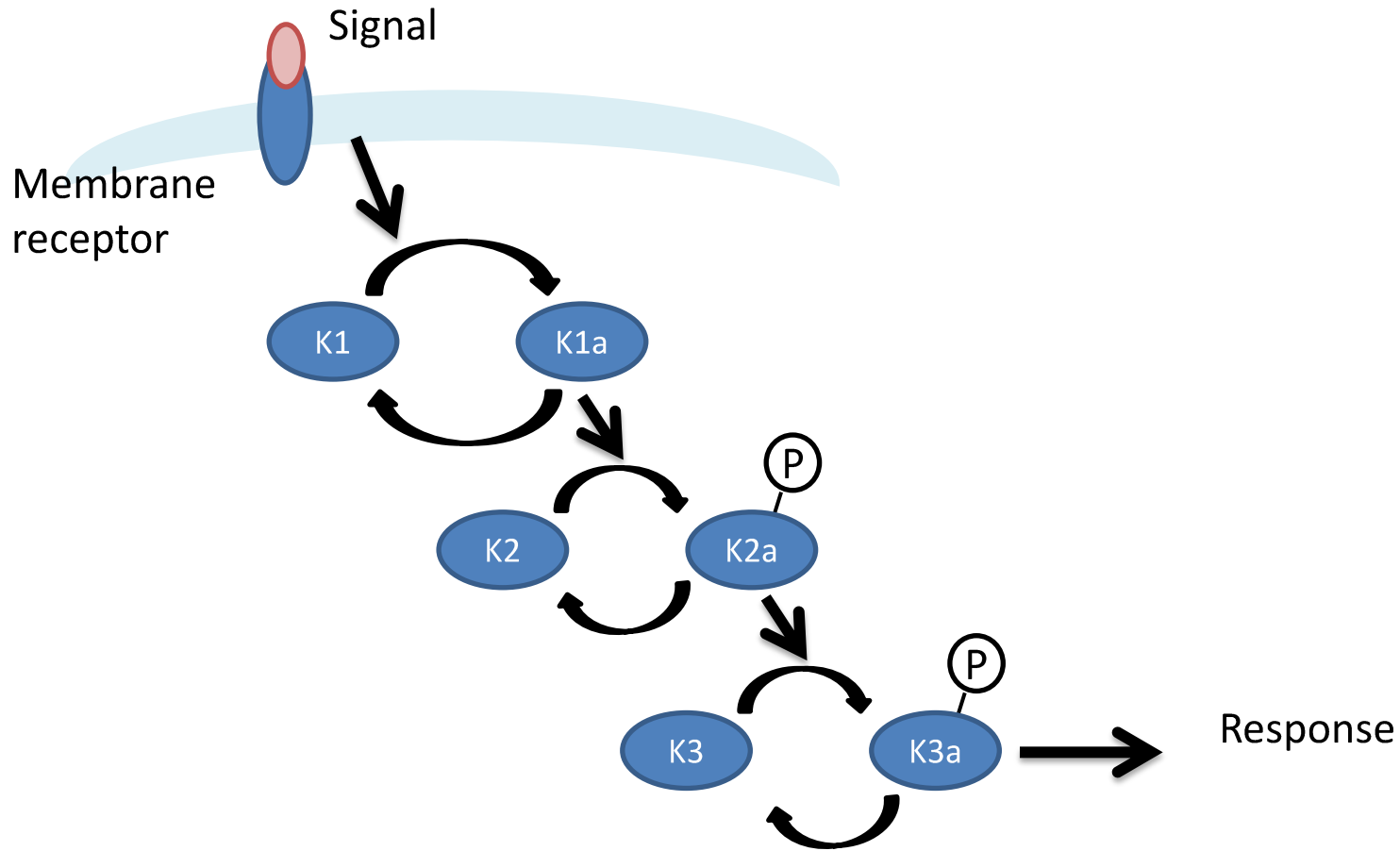
Partly explained by the presence of other amino-acids around the Serine (specificity “motifs”).

Example of kinase “motifs”

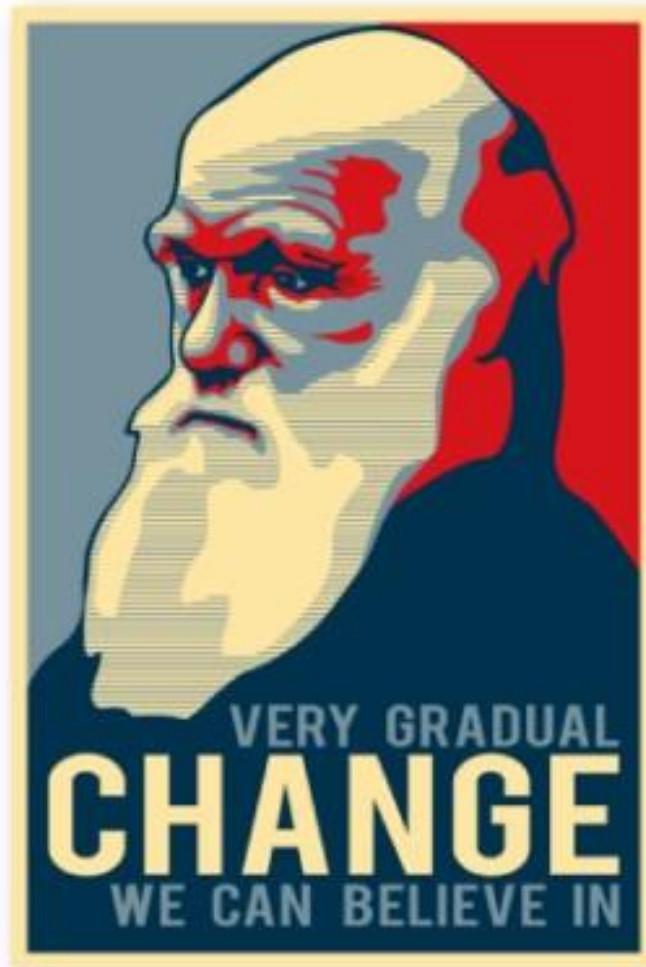
Cdk kinase - (K/R)S*PX(K/R)

These preferences are not the only specificity determinants . Other factors are important (co-localization, co-expression, interactions with common partners, etc)

Cascades of PTMs



How do species evolve ?
How are new functions created ?







Point mutations
Recombination
Duplications

Mutants

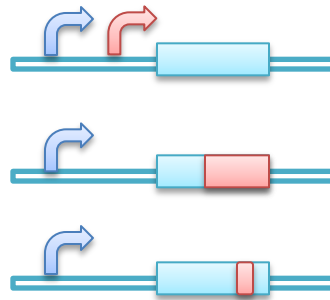




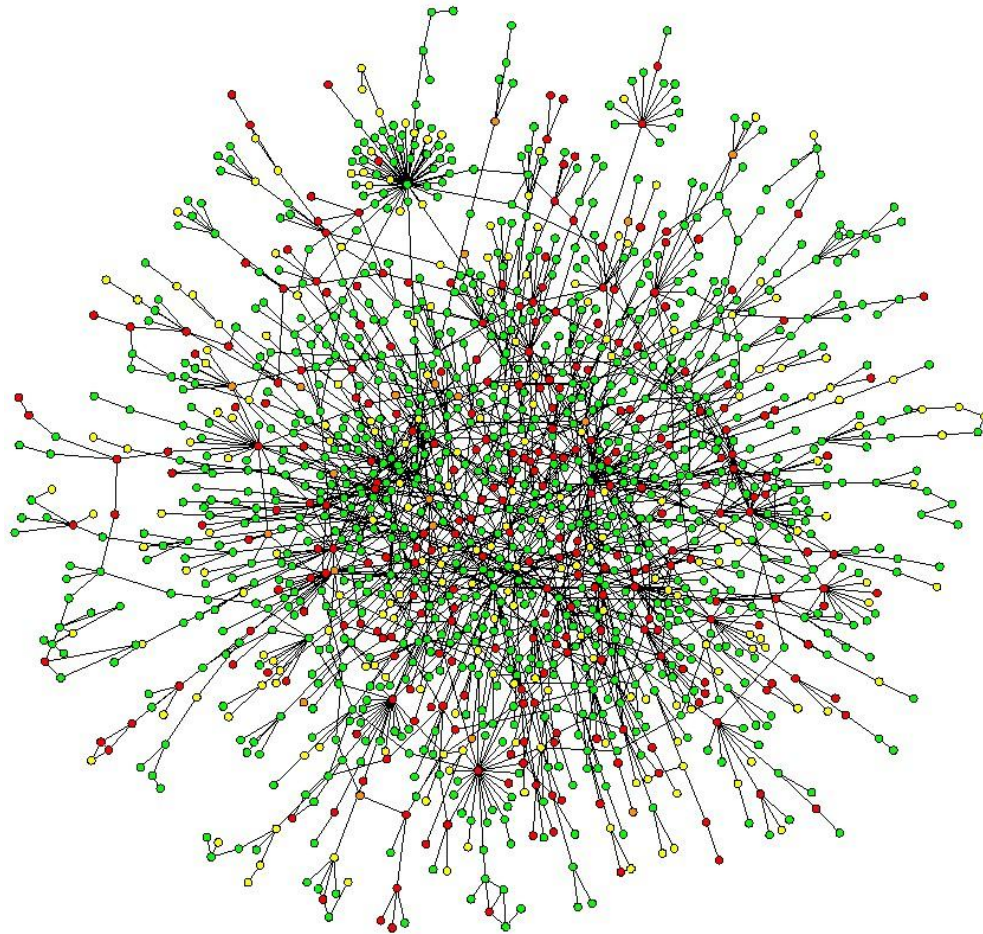
Mutants



Point mutations
Recombination
Duplications



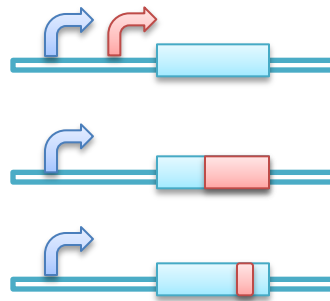
Changes in protein-protein,
protein-DNA interactions, etc



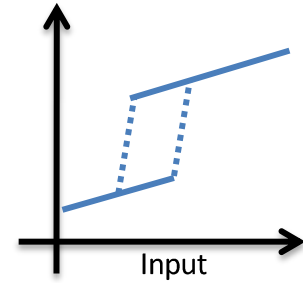
Cellular interaction networks



Point mutations
Recombination
Duplications



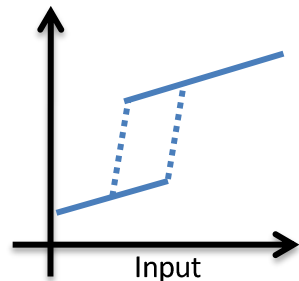
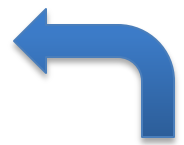
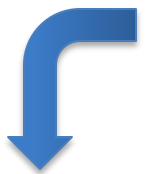
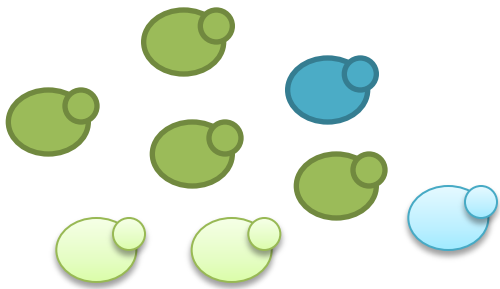
Changes in protein-protein,
protein-DNA interactions, etc



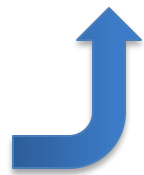
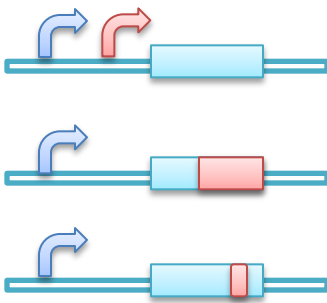
Phenotypic changes



Fitness differences



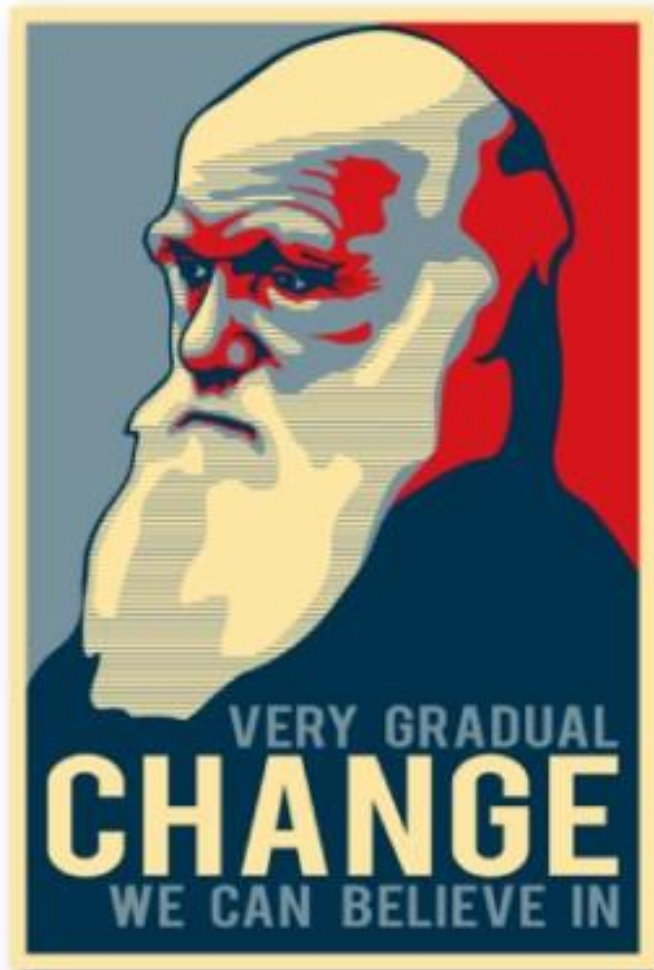
Point mutations
Recombination
Duplications



Changes in protein-protein,
protein-DNA interactions, etc

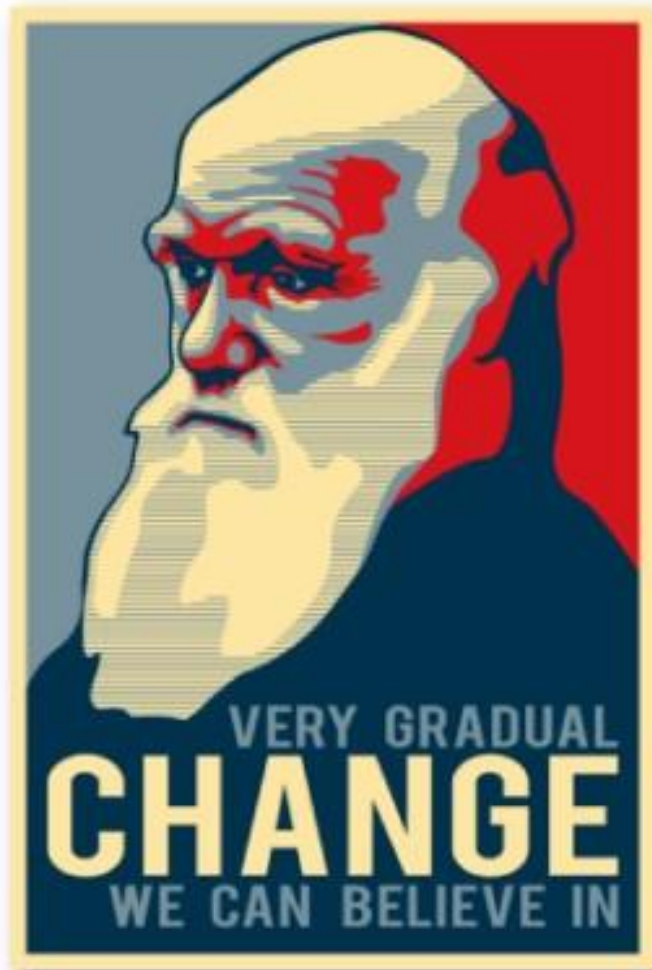
Phenotypic changes

Evolution of PTMs



What questions
can we ask ?

Evolution of PTMs

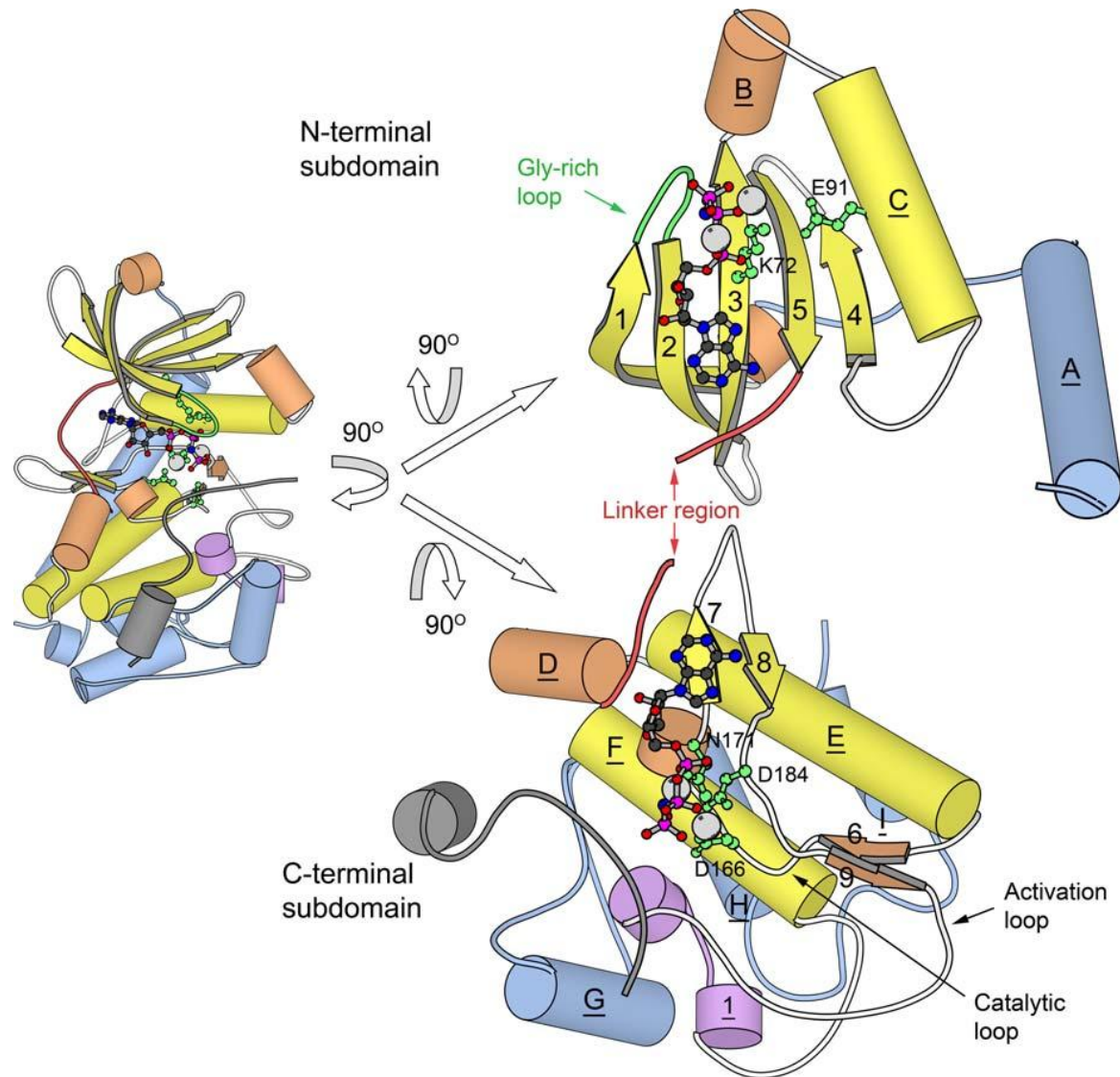


What questions can we ask ?

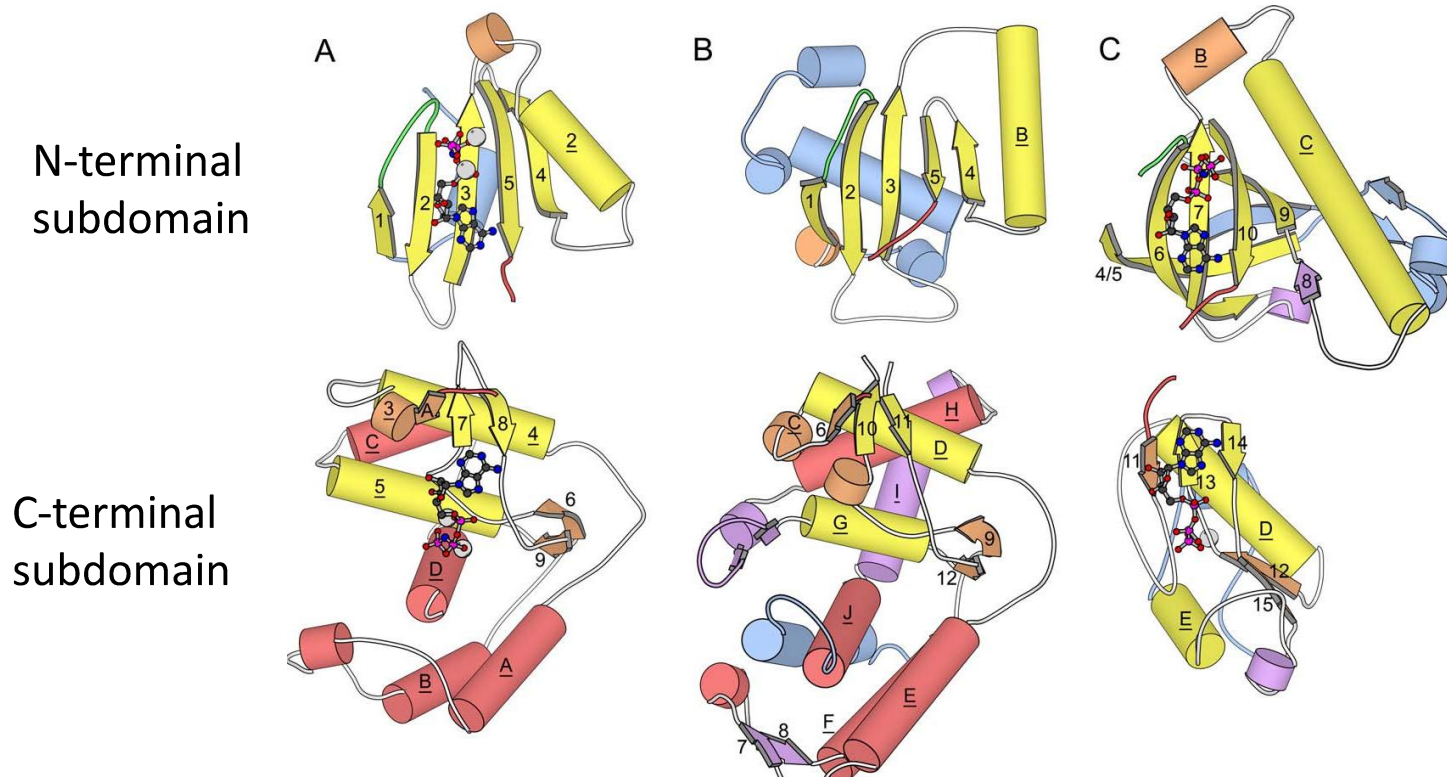
- What is the origin of the PTM ?
- How do the regulators change ?
- How do the targets change ?
- What changes have impact on fitness ?

I will focus on phosphorylation as the most well studied example.

Origin of Phosphorylation studied by structural homology

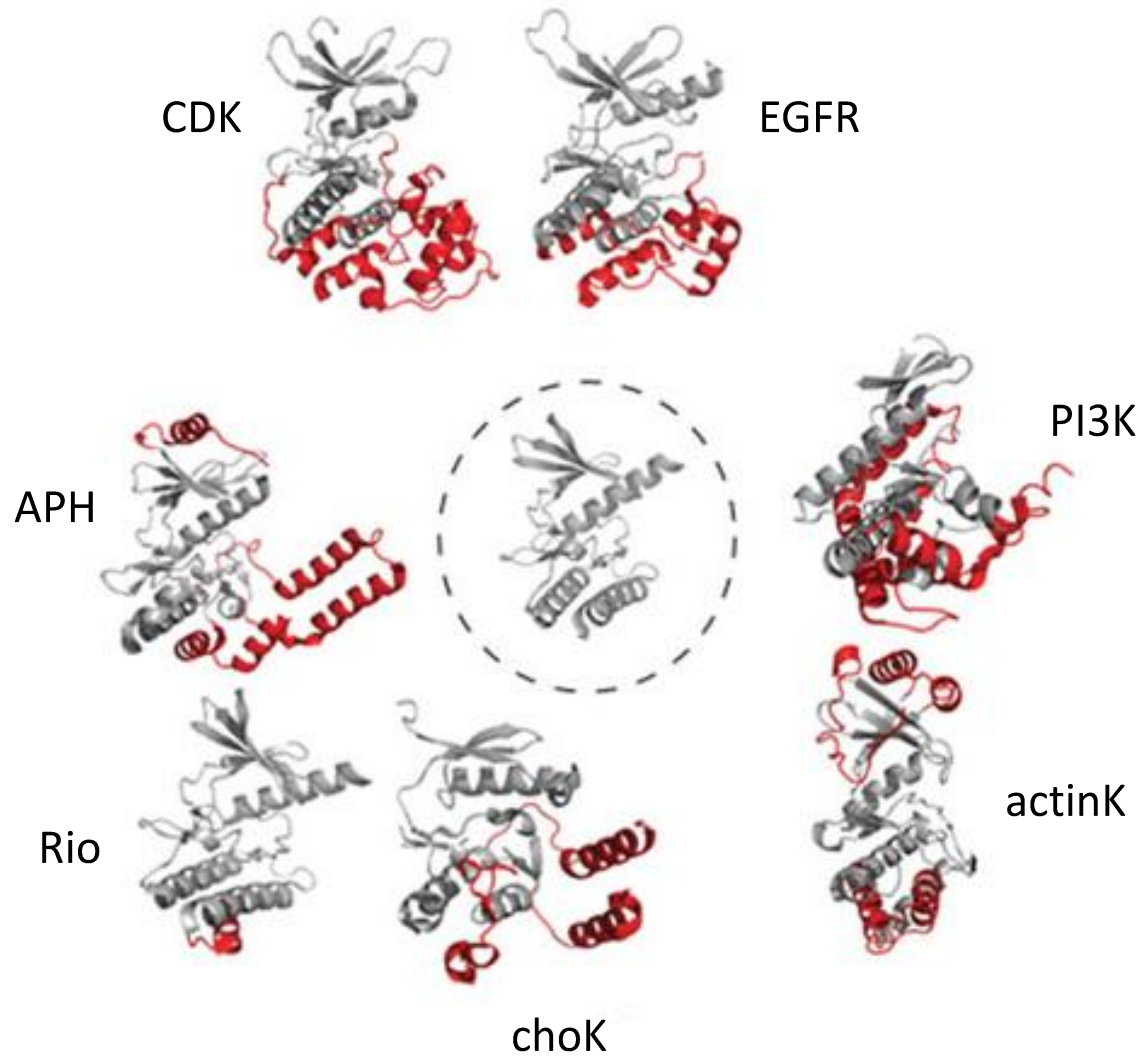


Origin of Phosphorylation studied by structural homology

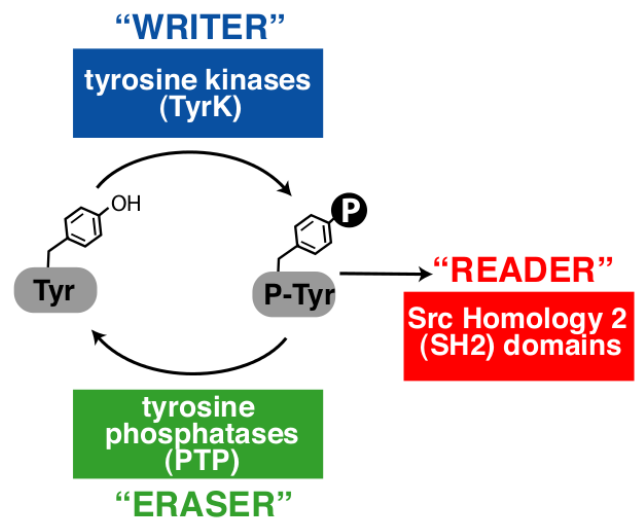
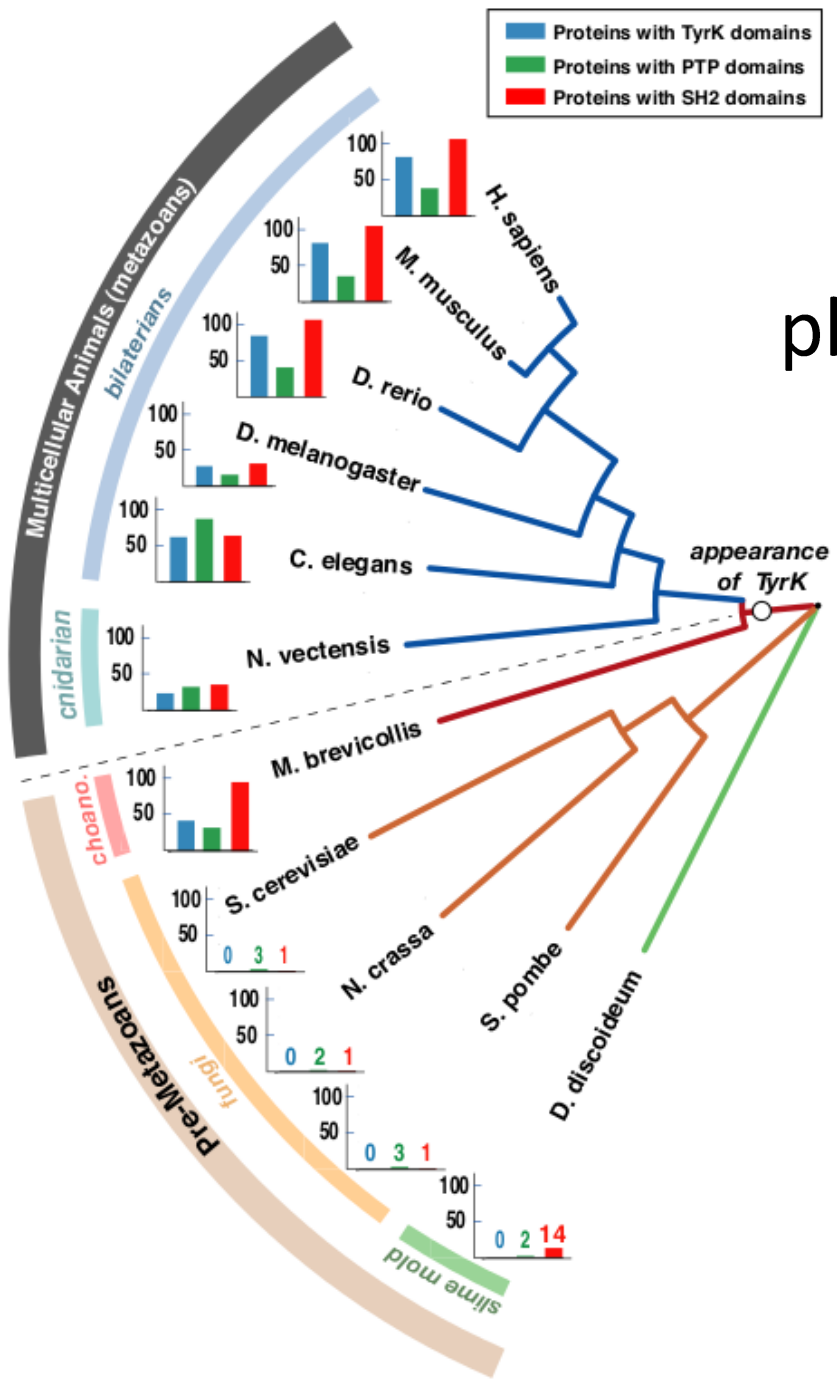


Protein kinases exist across all life but have been expanded significantly in eukaryotic species. There are kinases that phosphorylate small-molecules and lipids. It is hard to study the origin of these very ancient events

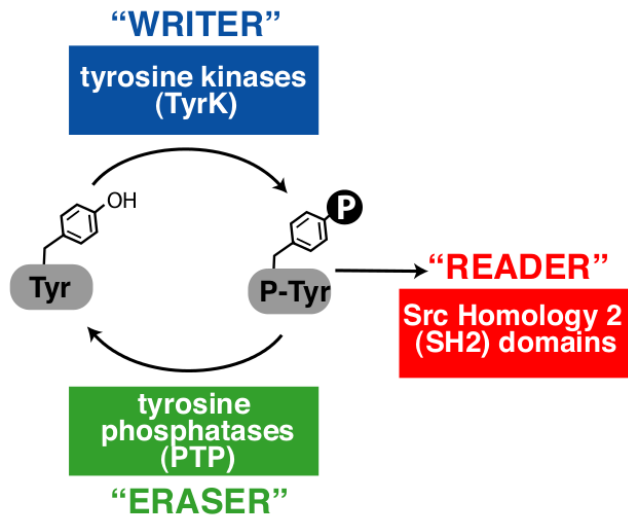
Diversification of protein kinases (through gene duplication and divergence)



Evolution of the tyrosine phosphorylation machinery

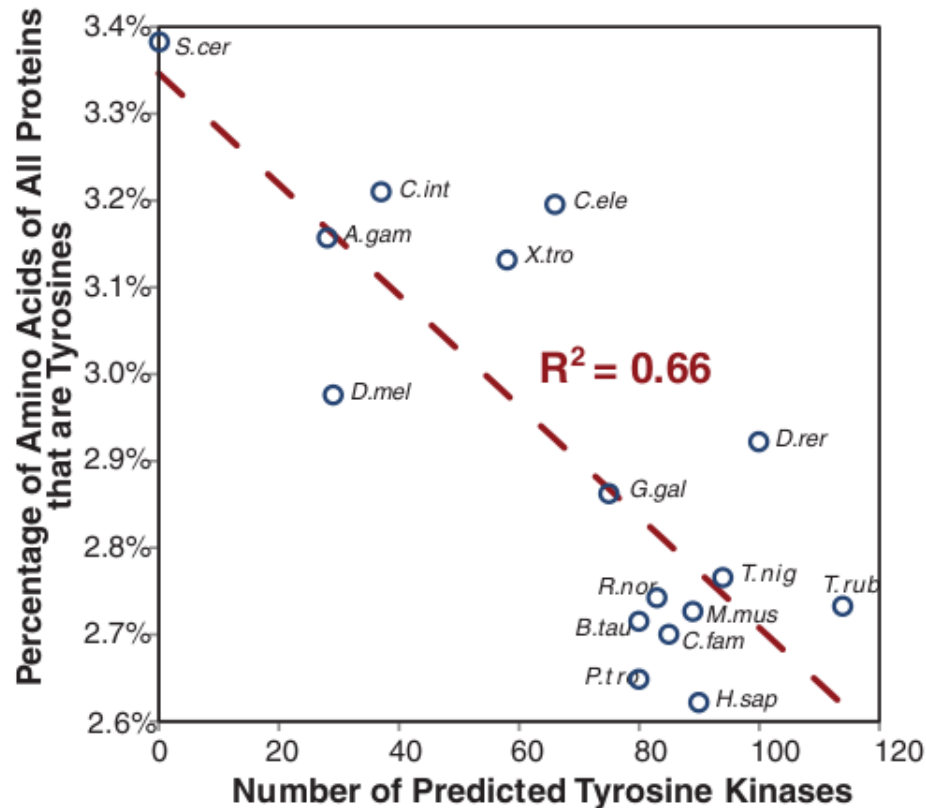


Costs of having a new modification



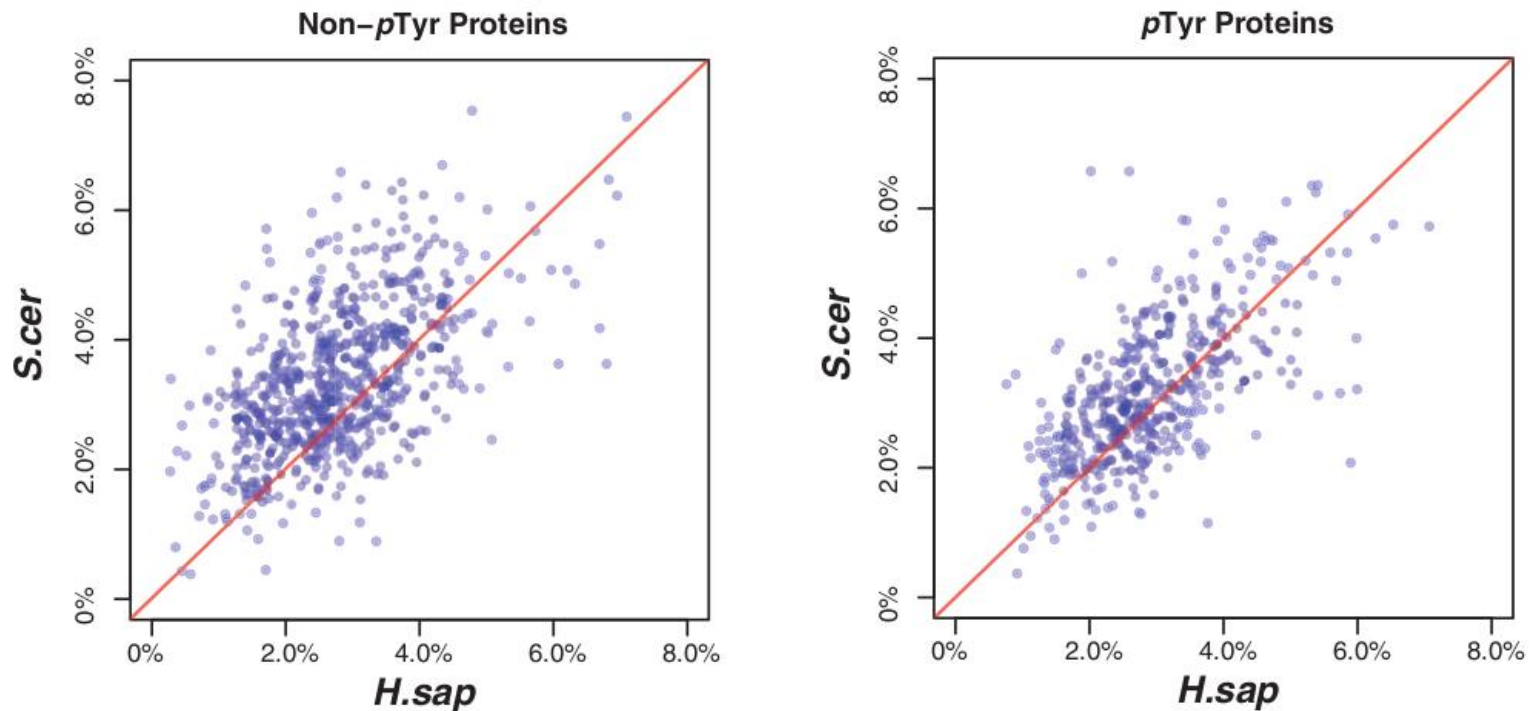
What sort of problems did the first tyrosine kinases create after they were accidentally invented by cells ?

Costs of having a new modification



Species that have more tyrosine kinases have proteins with lower percentage of tyrosines

Costs of having a new modification

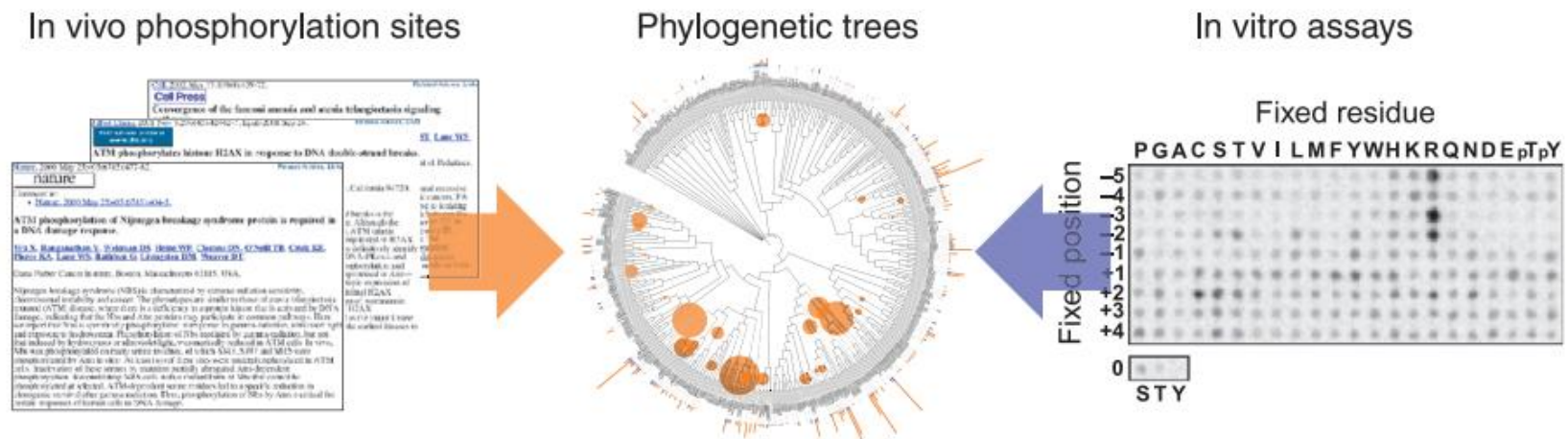


The reduction of tyrosines (in the yeast orthologs) is more noticeable in proteins that are today phospho-tyrosine proteins in human. Suggests: Some tyrosine phosphorylation was detrimental and there was selective pressure to mutate these to other amino-acids.

Predicting the specificity of the subfamilies

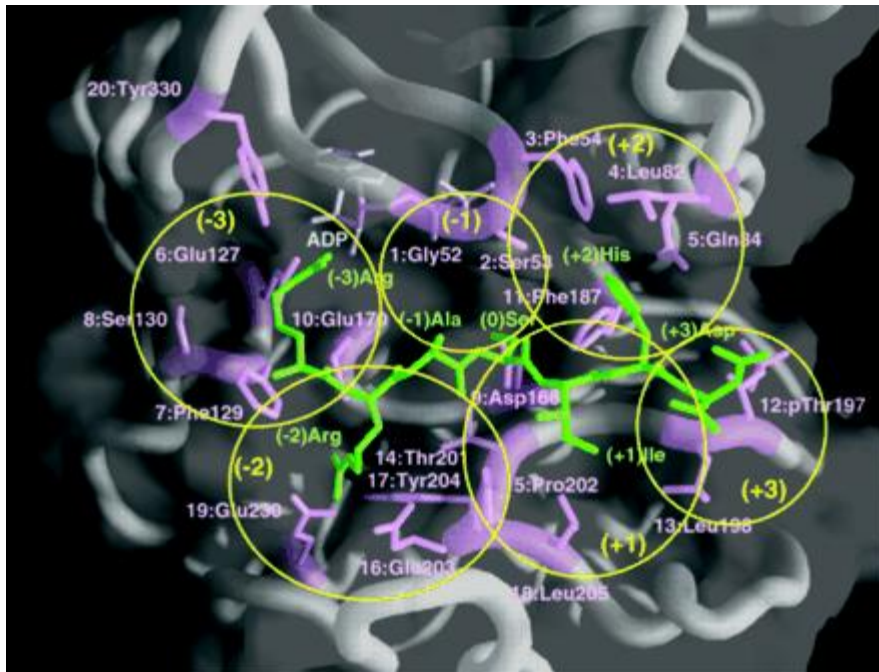
- Each kinase sub-family and (presumably) each kinase protein has a preference for specific targets. The specificity of the binding site can be predicted by two different approaches:
 - Machine learning applied to set of known target sites (e.g. Netphorest)
 - Structural analysis of the binding interfaces (e.g. Predikin)

Machine learning applied to set of known target sites (e.g. Netphorest)



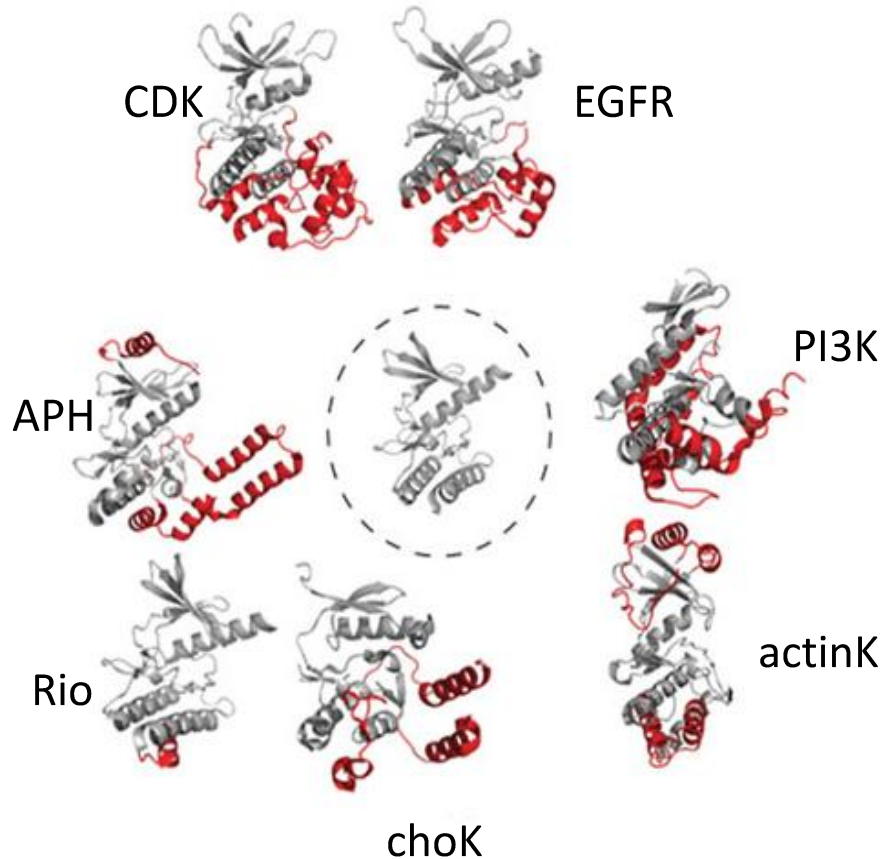
After reducing potential redundancies these were used as Training/test sets for data mining

Structural based prediction (e.g. Predikin)



- Compiled a set of structures of kinases bound to substrate sites.
- Extracted the residues that were in contact with the peptide.
- Correlated changes in kinase residues with preferences for the target “motif”
- Defined a set of rules that can be applied to different kinases based on sequence

Diversification of protein kinases (through gene duplication and divergence)



- These methods have so far not been applied to study the evolution of binding specificity
- They have been used to predict interactions between kinases and phosphosites and to study their evolution (in next section)

Recap – Evolution of regulators

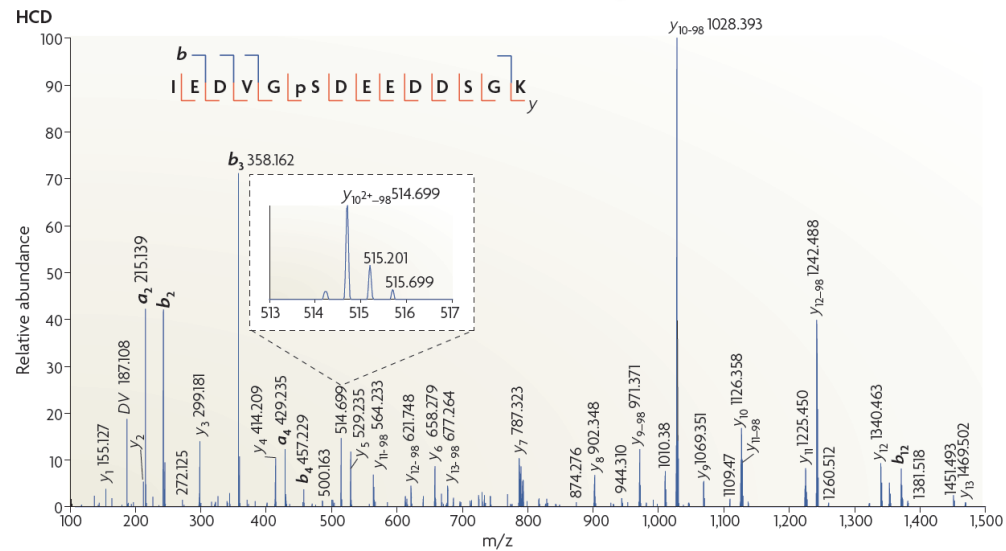
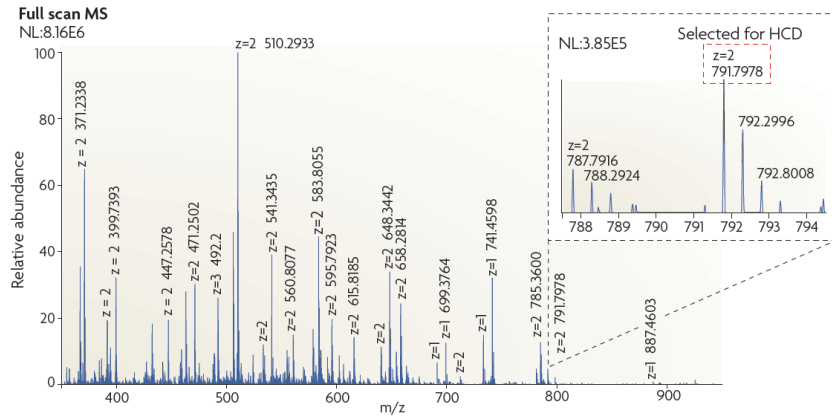
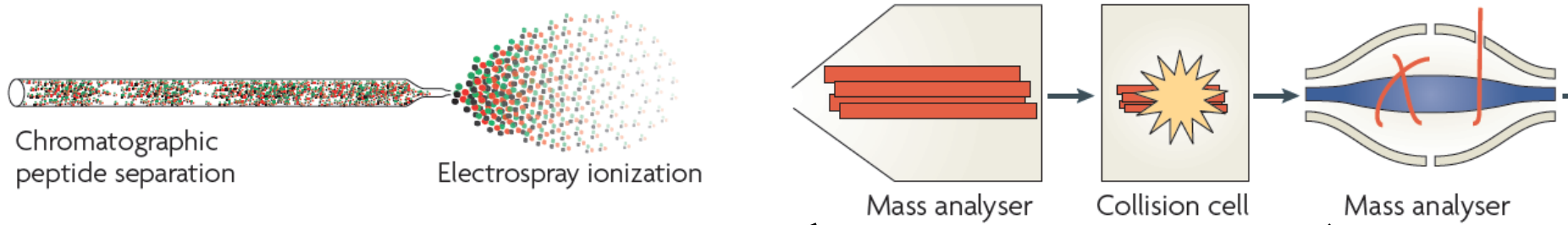
- Novel PTMs are created very rarely and their ancient origin is hard to study
- The usefulness of the PTM depends to some extent of having a set of 3 different types of proteins (writers, erasers and readers).
 - The origin of tyrosine phosphorylation appears to have coincided with a burst of duplication of the 3 types of domains
- Novel PTMs are likely to introduce pressures/costs since the proteome is not ready to have a large amount of the novel modification
- The PTM enzymes, once “invented” diverge by duplication and divergence into sub-families that have their own specificities
 - Kinase sub-families prefer to phosphorylate different amino-acid “motifs” around the phosphosite (ex. RS*PXX)
 - Kinases do not appear to be very specific and other types of information (co-localization, co-expression, etc) are important for recognition

Recap – Evolution of regulators (Computational Methods)

- The study of the very ancient evolutionary origins of PTMs depends on the analysis of structural homology.
- The duplication/divergence into subfamilies is commonly studied by training domain family models (ex PFAM database). These are used to search across multiple species.
- Structural bioinformatics or machine learning methods can be used to predict the specificity of kinases

Evolution of PTM sites and interactions

PTMs can be experimentally identified using mass-spectrometry



PTMs can be experimentally identified using mass-spectrometry

1. Sample preparation
 1. Extraction
 2. Purification
 3. Protein digestion (**peptides**)
 4. Enrichment for PTM (antibody or binding column)
2. Tandem Mass-spec (MS/MS)
 1. Chromatographic peptide separation
 2. Ionization
 3. Peptide ions mass/charge (m/z) determination
 4. Peptide fragmentation
 5. Fragment ions m/z determination
3. Computational analysis
 1. Compare experimental spectra with predicted spectra
 2. Determine false discovery rate

Cell

Resource

36,000 phosphorylation sites

A Tissue-Specific Atlas of Mouse Protein Phosphorylation and Expression

Edward L. Huttlin,^{1,4} Mark P. Jedrychowski,^{1,4} Joshua E. Elias,^{1,4,5} Tapasree Goswami,^{1,4} Ramin Rad,² Sean A. Beausoleil,^{1,6} Judit Villén,^{1,7} Wilhelm Haas,¹ Mathew E. Sowa,^{1,3,6} and Steven P. Gygi^{1,2,*}

Molecular Cell

Resource

19,000 ubiquitylation sites

Cell
PRESS

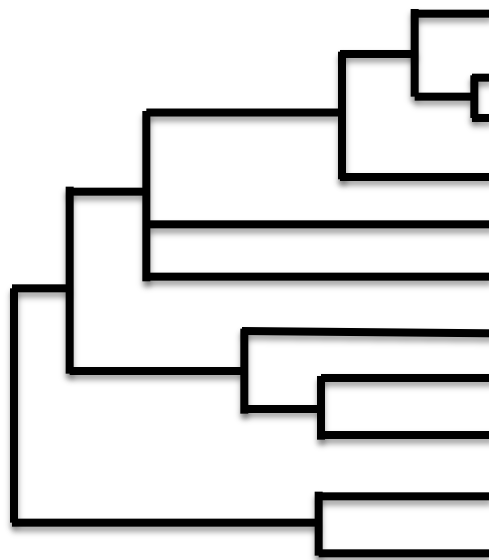
Systematic and Quantitative Assessment of the Ubiquitin-Modified Proteome

Woong Kim,^{1,4} Eric J. Bennett,^{1,2,4,5} Edward L. Huttlin,¹ Ailan Guo,³ Jing Li,³ Anthony Possemato,³ Mathew E. Sowa,^{1,2} Ramin Rad,¹ John Rush,³ Michael J. Comb,³ J. Wade Harper,^{1,2,*} and Steven P. Gygi^{2,*}

The enrichment methods and sensitivity of the MS machines has improved tremendously in the past 5 years.

Known eukaryotic PTM sites

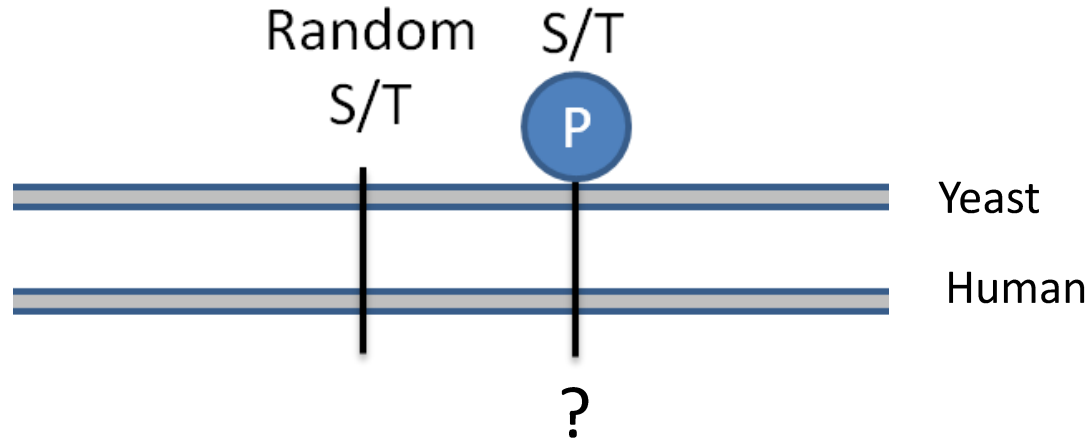
(curated from MS studies)



Species	Phospho sites	Acetylation	Ubiquitylation
<i>H.sapiens</i>	45943	8042	22000
<i>M.musculus</i>	32885	3384	
<i>R.norvegicus</i>	2702		
<i>X.laevis</i>	629		
<i>C.elegans</i>	8433		
<i>D.melanogaster</i>	22307	1707	
<i>S.pombe</i>	2989		
<i>S.cerevisiae</i>	22536		5000
<i>C.albicans</i>	3773		
<i>A.thaliana</i>	5336		
<i>O.sativa</i>	3666		

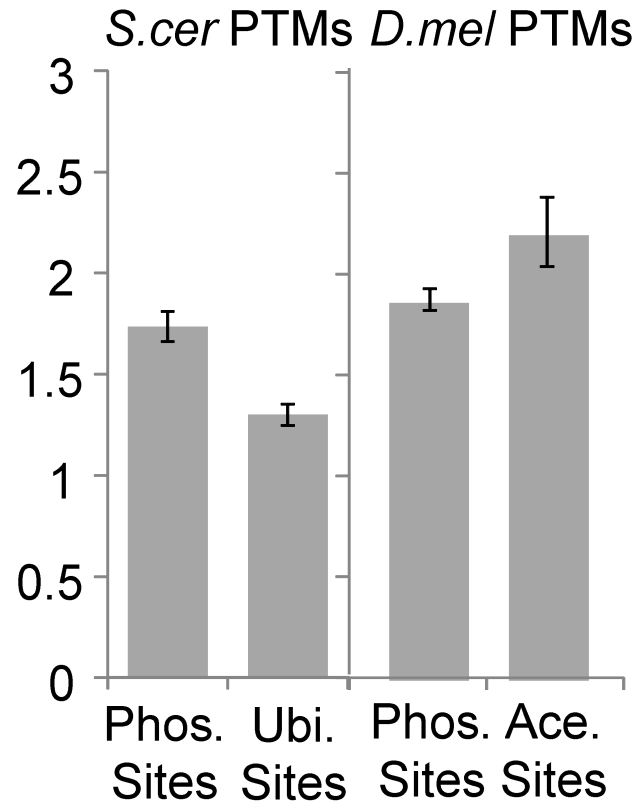
> 150.000 phosphorylation sites
 ~ 200.000 PTMs

Measuring conservation



- Conservation of “State” (ex. phosphorylated): The same peptide in the human protein is also phosphorylated
- Compared with random sampling of acceptor residues in the same proteins - ratio of conservation over random.

Conservation of post-translational modifications



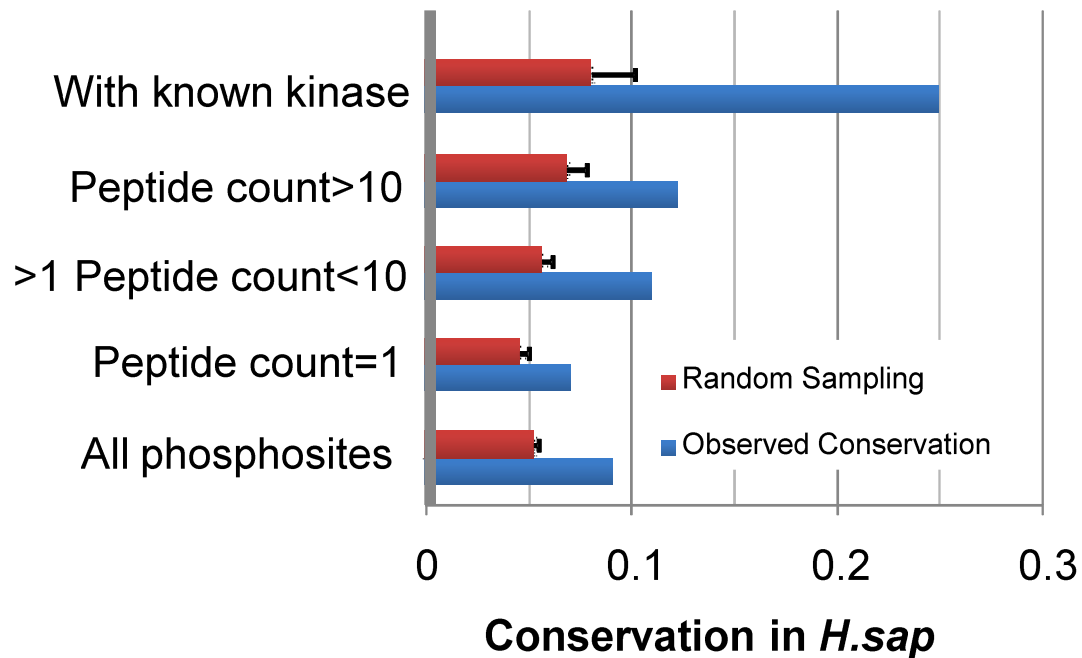
All 3 modifications are poorly evolutionarily constrained when compared to random sample of acceptor residues

If there is a difference:
AceK > phospho > ubi

Suggest that ubiquitylation is less specific than the other 2.

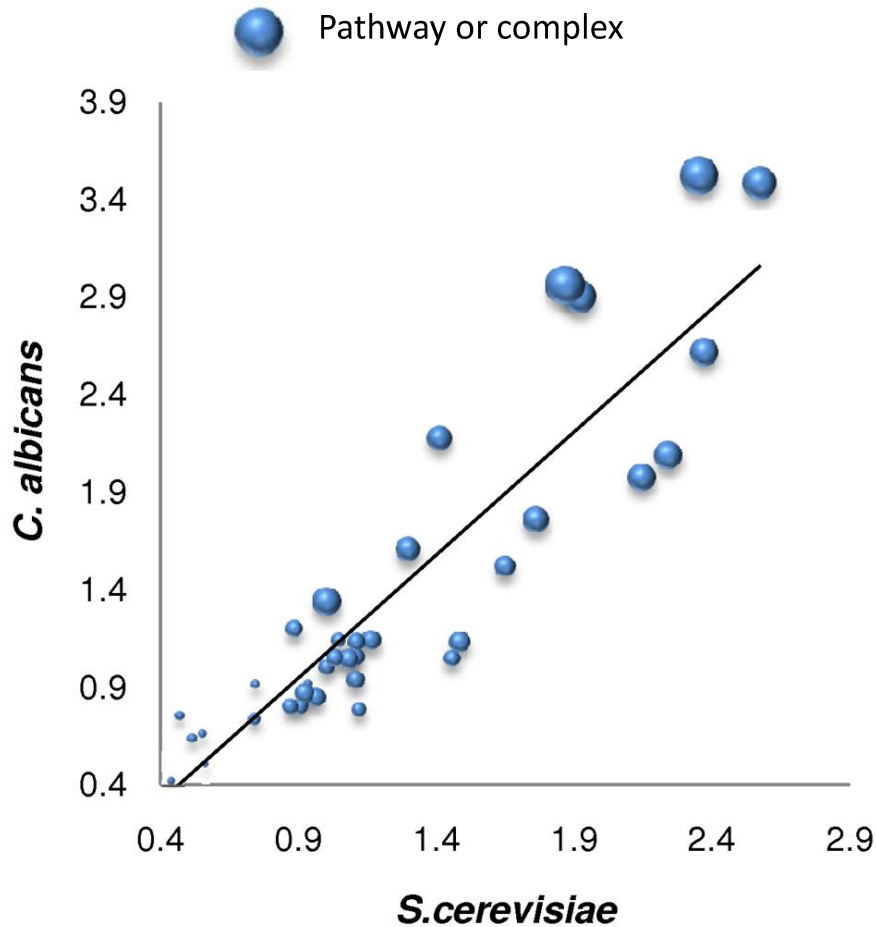
Phosphosites with known function are more likely conserved

S. cerevisiae phosphosites



It is possible that many sites have no biological function

Phosphorylation of functional 'modules' is well conserved



Average number of phosphosites per protein
(normalized for proteome coverage)

It is also possible that different sites have the same function in different species.

Phosphosite conservation and protein abundance

- Some fraction of phosphosites are not functional
- Protein abundances span different orders of magnitude

Question: Are phosphosites more conserved in proteins of high or low abundance ?

Evolution of kinase-substrate interactions

In order to study the evolution of kinase-substrate interactions. We need to predict what kinases are responsible for the phosphorylation.

How would you predict kinase-site interactions ?

Prediction of kinase-site interactions



+

Co-expression
Co-localization
Scaffolding
...

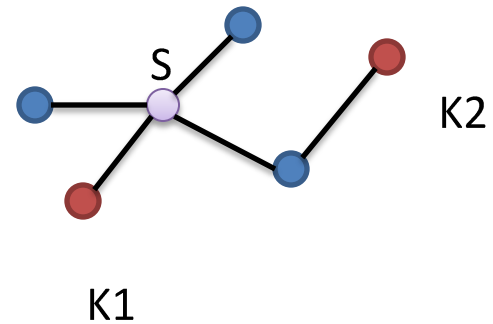
In vitro

In vivo

Prediction of kinase-site interactions



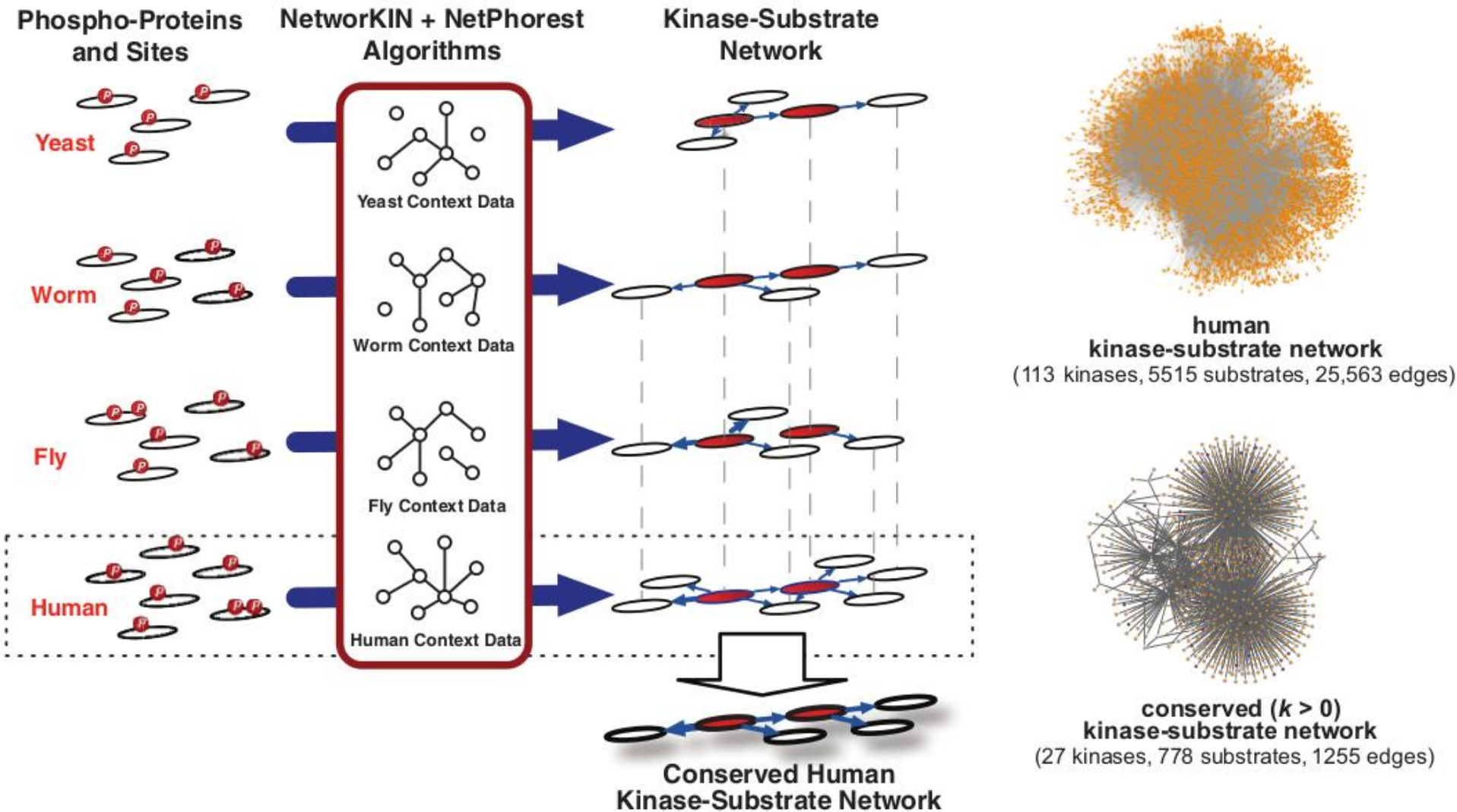
+



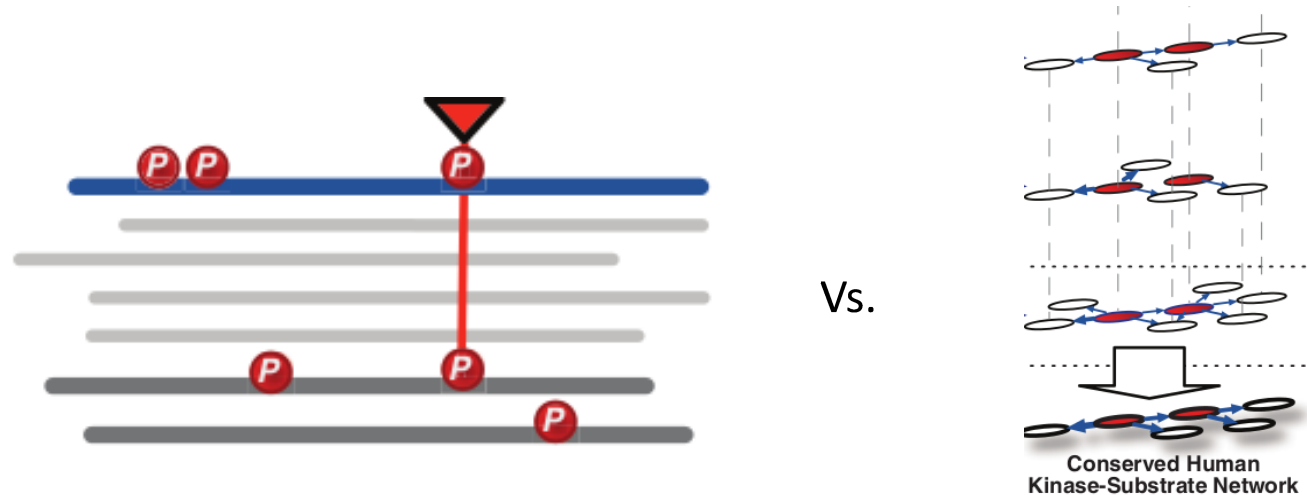
Models
from
Netphorest

Functional
interactions
from
STRING

Evolution of kinase-substrate interactions

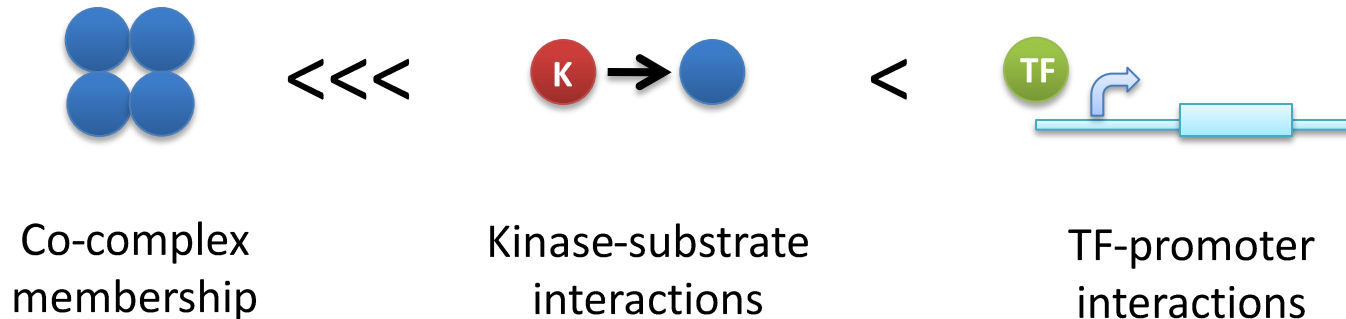


Evolution of kinase-substrate interactions



- Only a small fraction of the phosphorylation sites are conserved in the same position of the alignment
- A significant fraction of the kinase-substrate interactions are conserved but with predicted changes in phosphosite position.
- The network of conserved kinase-substrate interactions is enriched in proteins that when mutated can result in disease.

Evolution of different types of physical interactions



The specificity of interaction determines the likelihood that a new interaction can be created by random variation. Lower specificity correlates with higher rate of divergence of interactions.

TF-gene Interactions diverge very quickly during evolution

S. cerevisiae *K. lactis* *C. albicans*

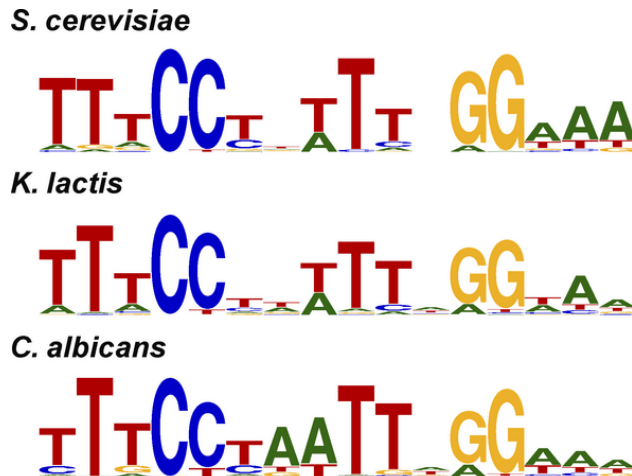
Fraction of genes mapped from species A bound in species B

<i>S. cerevisiae</i>		42%	22%
<i>K. lactis</i>	16%		19%
<i>C. albicans</i>	7%	22%	

ChIP-Chip of Mcm1 (TF) in 3 different species (100 to 400 My).

Very poor conservation (~20%) of interactions although binding site specificity appears conserved.

We don't know how much this change impacts on phenotypes.



Protein Complexes are largely conserved

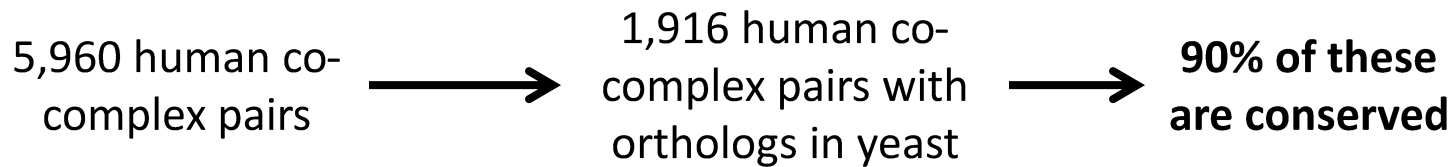
OPEN ACCESS Freely available online

PLOS COMPUTATIONAL BIOLOGY

Protein Complex Evolution Does Not Involve Extensive Network Rewiring

2008

Teunis J. P. van Dam^{1,2}, Berend Snel^{1,3*}

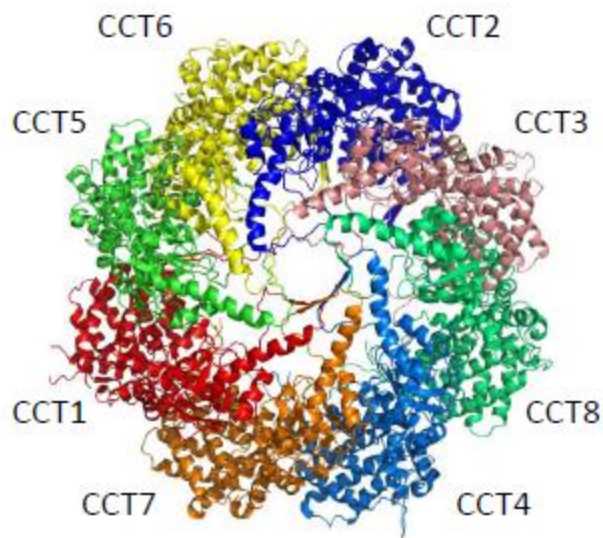


Protein complex evolution by duplication-divergence

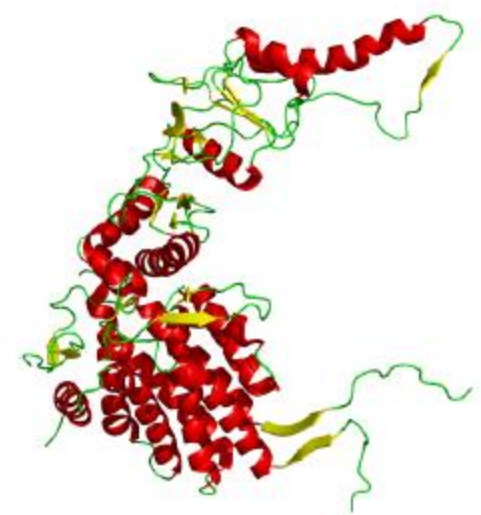
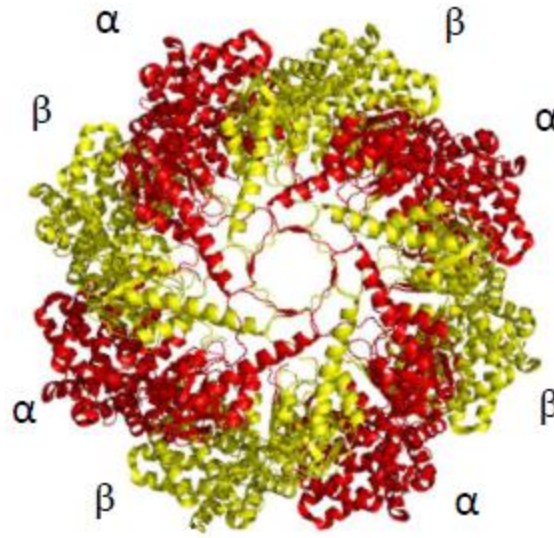
S.cerevisiae CCT/TRiC complex

Thermococcus strain KS-1 chaperonin

CCT1 subunit (hsp60 fold)

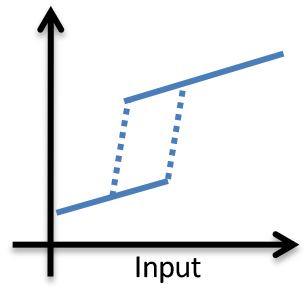
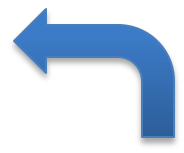
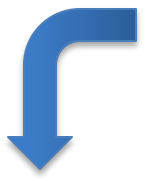
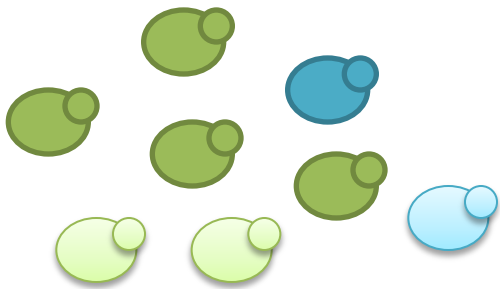


Top-view of complex

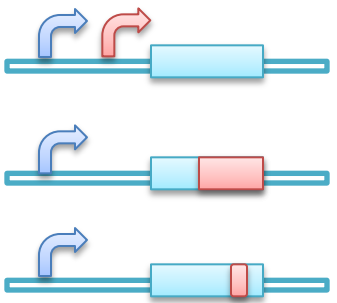
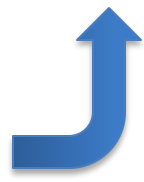


Side-view of subunit

Fitness differences



Point mutations
Recombination
Duplications



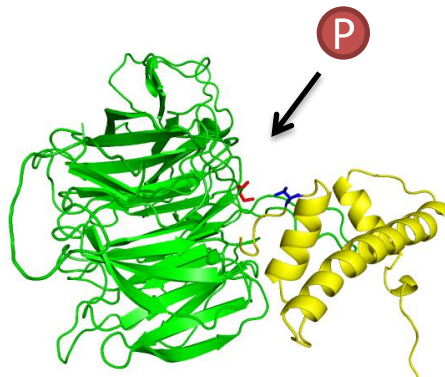
Changes in protein-protein,
protein-DNA interactions, etc

Systematic functional annotation of PTMs

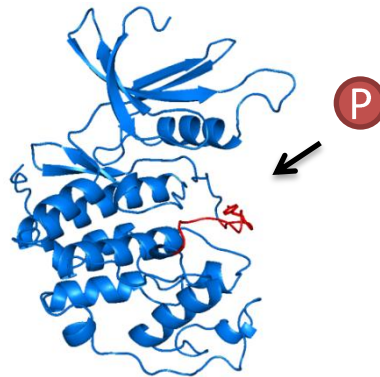
- PTMs diverge very quickly and it is likely that a significant number maybe have no function.
- To study the impact on fitness of these changes and for practical functional studies we need computational ways to predict the ones with function
 - Conservation of site and prediction of kinase (as described above)
 - **Assign a mechanist role for PTM**

Systematic functional annotation of PTMs

Globular regions

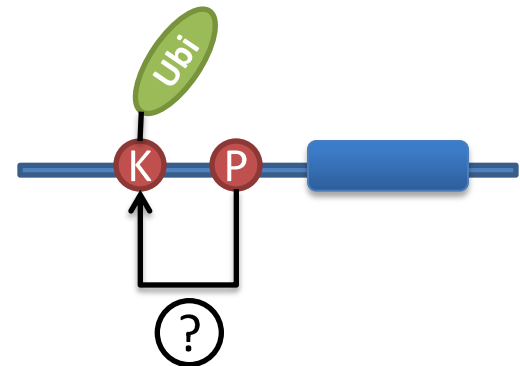


Regulating
interfaces



Regulating domain
activity (eg. allostery)

Unstructured regions



Phosphorylation switches

~25%

~45%

Phosphorylation

AceK/Ubi

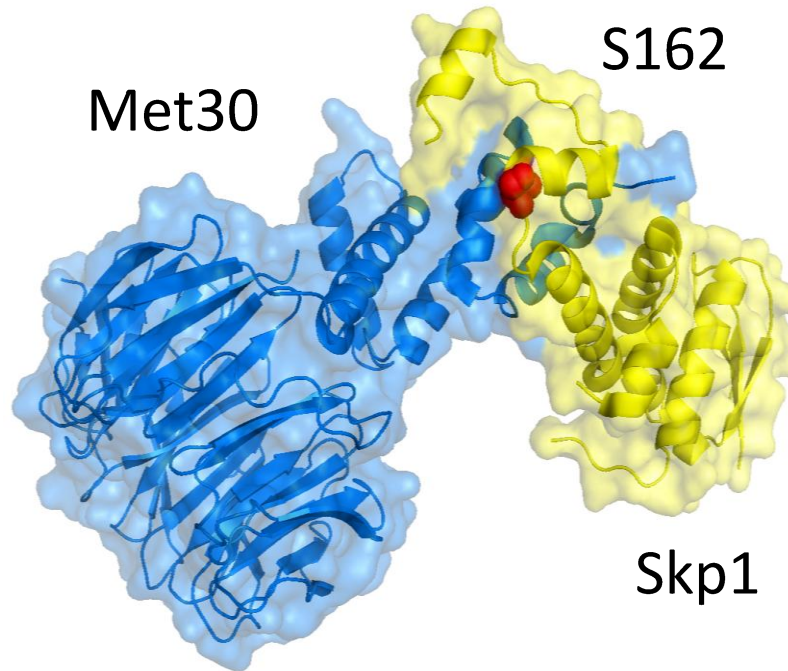
~75%

~55%

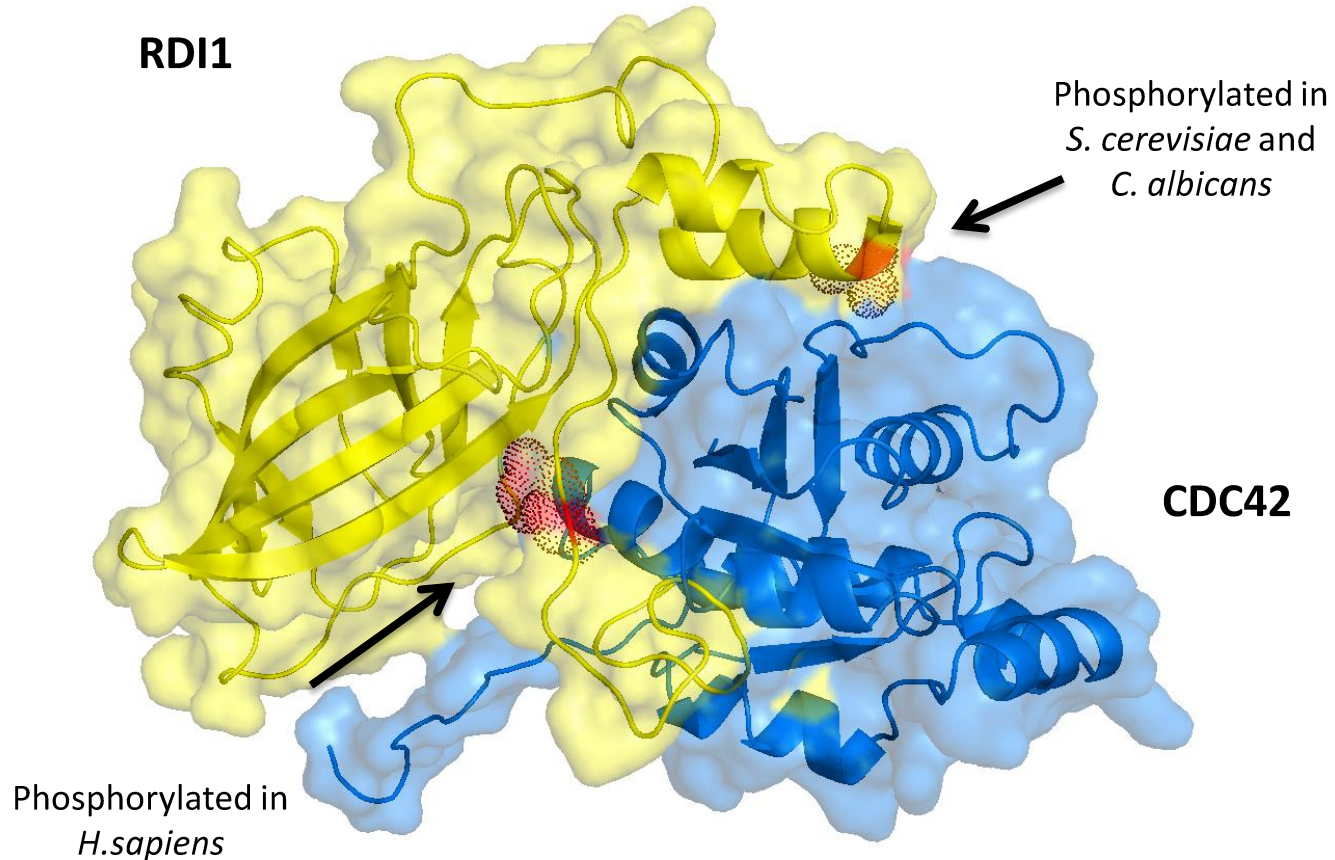
PTMs at interface residues

Interface residues were predicted based on xray, homology models, docking solutions and domain-domain contacts from the 3DID database

Example:
Skp1:Met30



Conservation of interface PTMs vs. conservation of function



Phosphorylated residues are 20 Amino-acids apart in sequence

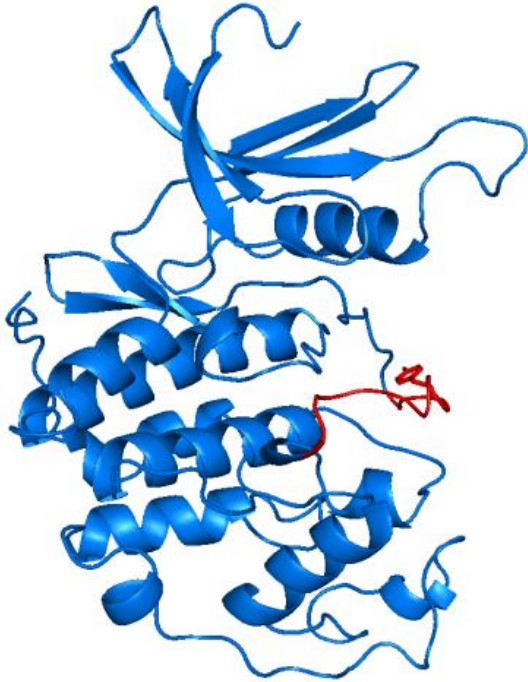
Conservation of function vs. conservation of phosphosites

Phosphosite	In <i>S. cerevisiae</i>	Conserved in another species by aln (+/- 5 pos)	Fraction conserved
All	22536	3585	0.16
In interface	401	142	0.35

	In <i>S. cerevisiae</i>	Conserved in another species	Fraction conserved
Interfaces with pSites (xray)	6	4	0.66
Interfaces with pSites (docking)	261	131	0.50
Interfaces with pSites (comparative modeling)	107	81	0.76

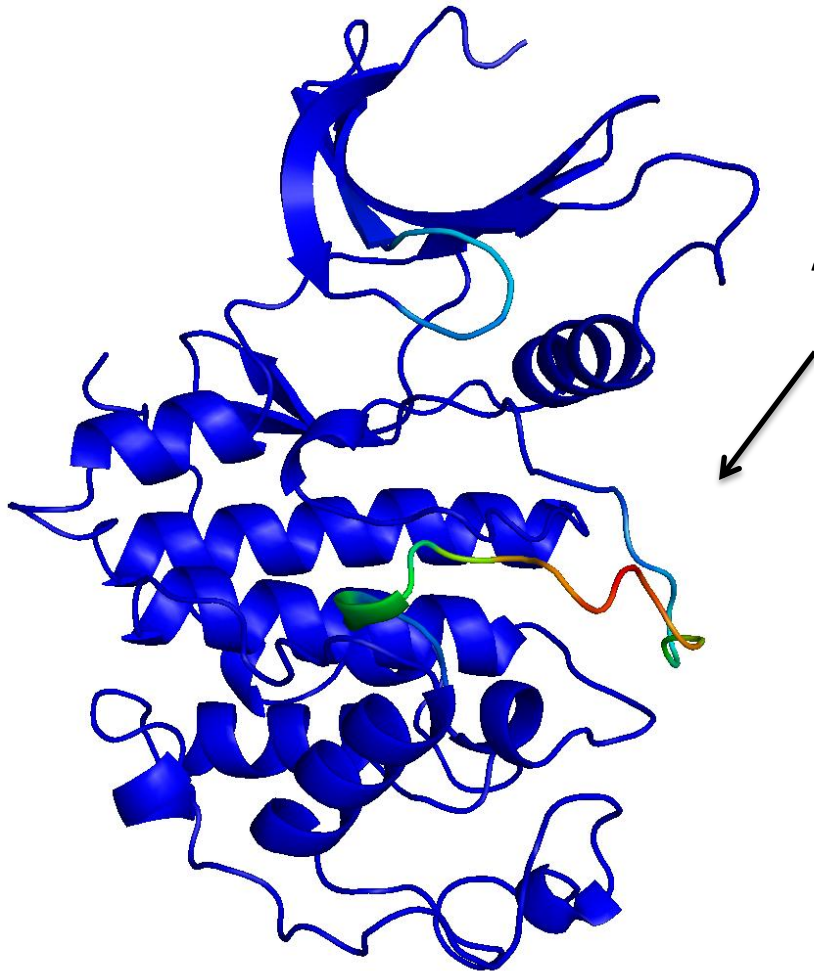
Positional conservation under-predicts conservation of function
It is possible to change a PTM site without changing its function

Regulation of domain activity

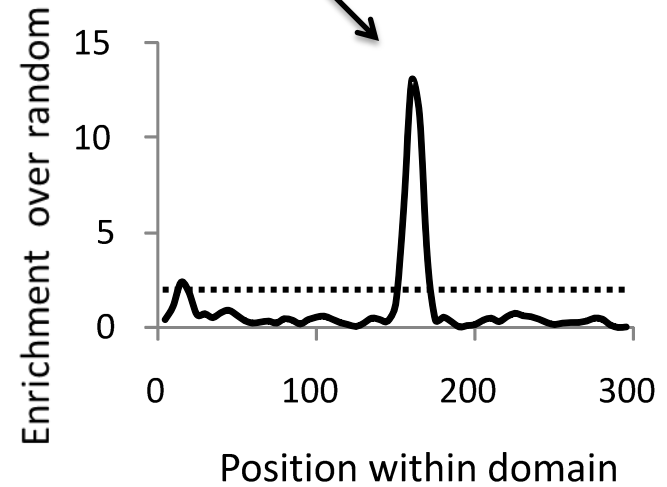


1. Pool all phosphosites (or other PTMs) and map them to representative structures of domain families
2. Look for sequence / structural enrichment (indicative of conservation)

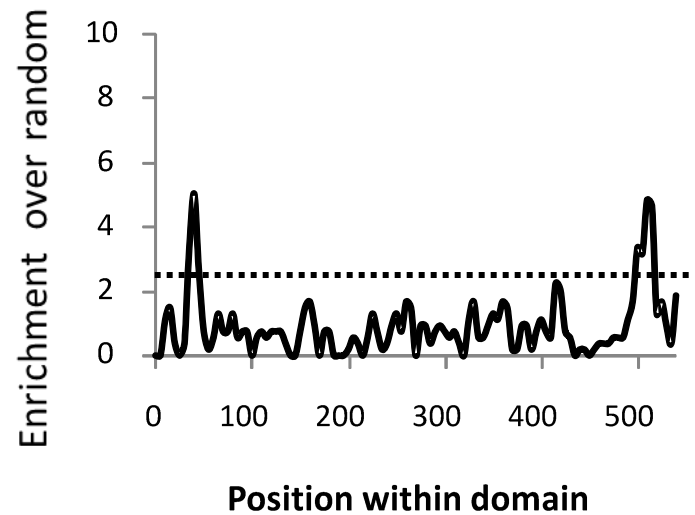
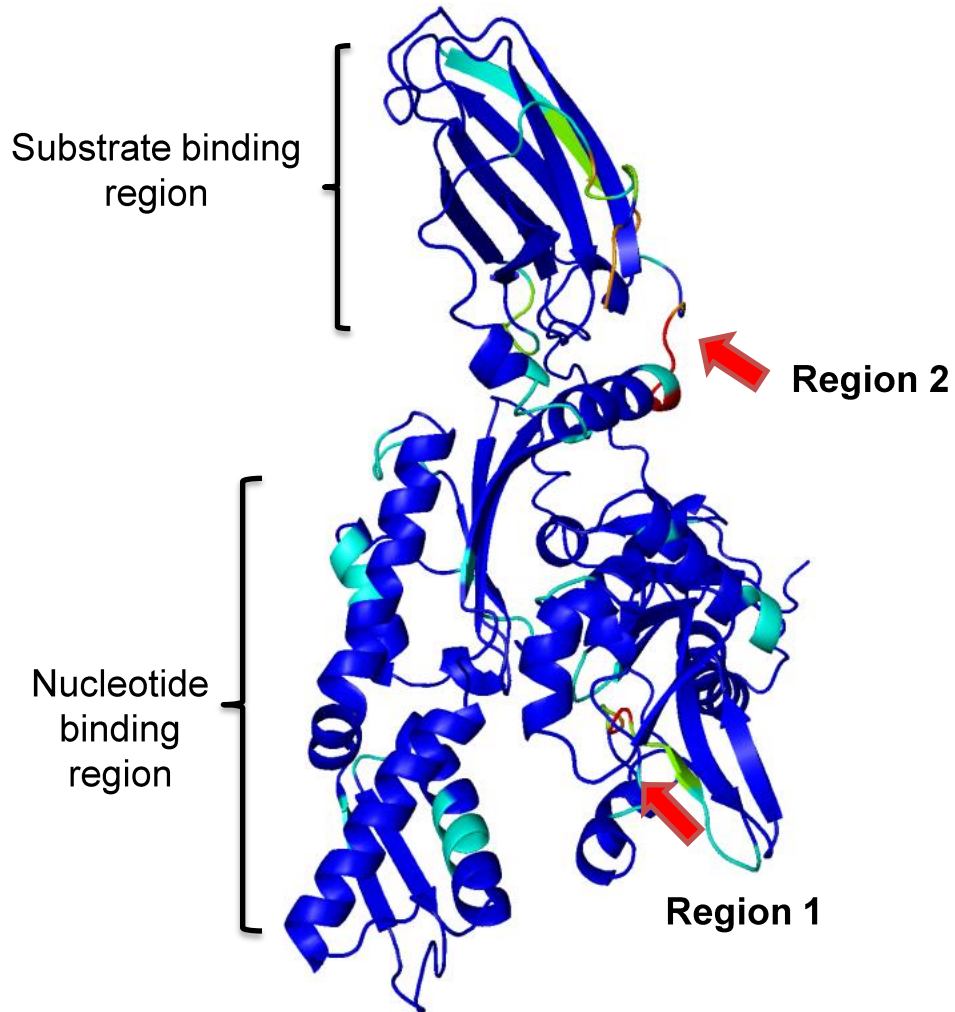
Protein kinase domain phosphorylation propensity



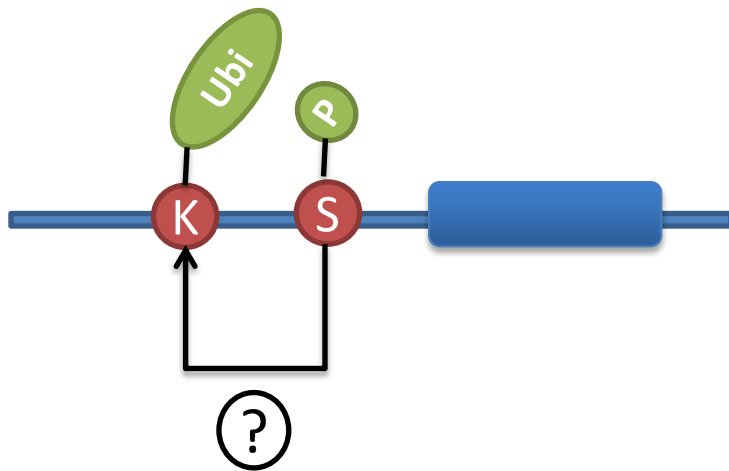
Activation loop (50% of psites)
Regulatory 'hot-spot'



HSP70 domain Phosphorylation propensity



Co-regulation between different PTMs



We showed that:

- Different PTMs tend to cluster within proteins.
- Phosphosites near lysine modifications are more conserved than average phosphosites.

- The most promising but also the most challenging aspect of this analysis.
- New mass-spec and computational methods are needed (ex. dynamics after perturbation).

Recap - Evolution of targets

- PTM sites diverge very quickly at a rate that might be similar to the divergence of TF binding sites.
- Enzyme-protein interactions are more likely conserved than individual PTM positions
- PTMs with identifiable function shown an above average conservation
 - It is possible that a significant fraction have no function and it is possible that some sites change position without changing the function.
- Conservation of sites and interactions can be used to find important PTMs

Recap - Evolution of targets (Computational Methods)

- Standard comparative genomics has been used to study the evolution of PTM positions
- Machine learning methods have been used to predict kinase-protein interactions to study their evolution
- Structural bioinformatics has been useful to identify sites that are more likely functionally important.
- A lot of work still needs to be done in ranking sites/interactions according to functional importance as well as assigning the most likely mechanistic work.