

## Use expression Python

## Regular Expressions

a b	"a" or "b"	
?	Zero or one	
+	One or more	
*	Zero or more	
*?	Zero or more, but stop after 1, matches	
{N}	exactly N occurrences	{, 3} up to 3
{N, n}	from N to n occurrences	

[A-Z] uppercase from "A" to "Z"  
[0-9] any number  
[a-zA-Z] any character in the list  
[^\a-zA-Z] any but not a character in the list  
. any character, except newline

- ^ start of a line
- \$ end of a line
- \s whitespace
- \S non-whitespace
- \w word: alphabet + underscore
- \W non-word; i.e. special etc.

1d digit  $\sim [0-9]$   
 2d num-digit

$$w \sim [A-Za-z0-9\_]$$

(...) capture groups (3 <sup>are</sup> chunks)

input re

S = "Biology = Sx fun!"

```
re.search("\d+", s) # return if match
# <re.Match object; span=(0, 1), match='5'>
```

# Search into genes/genomic coordinates

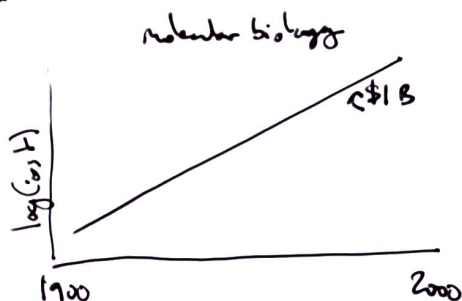
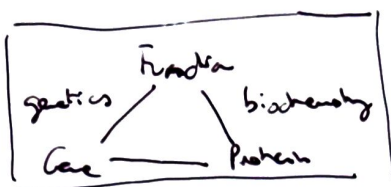
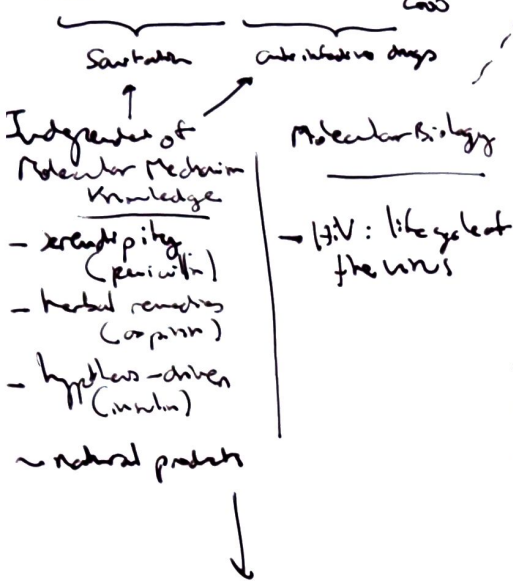
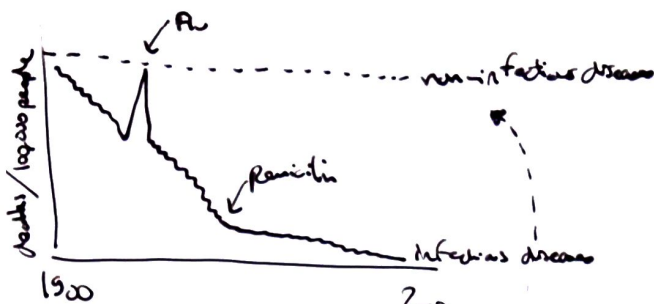
```
row = "50x2_chr3_34548926_34551382_+"

```

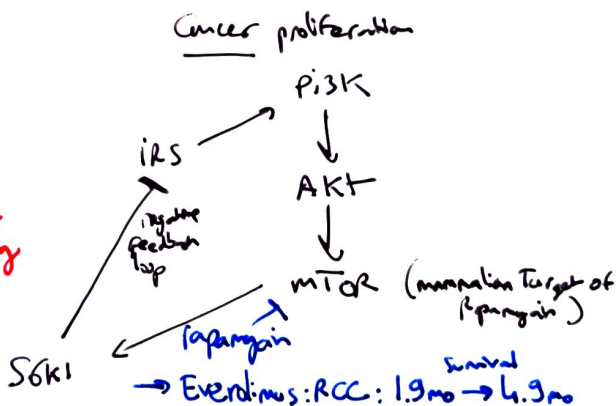
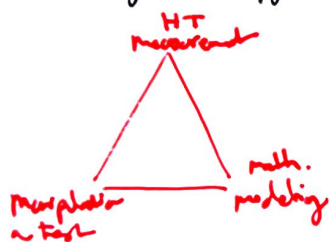
re. search("\w+|s{2}|chr|\w{,23}|s{2}(\d+|s{23}){2}[\+|-],\w+)

# Also working the following regex, to minimize potential for unexpected matches

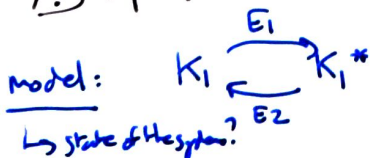
# \S+ \b+ chr \S+ \b+ \d+ \b+ \d+ \b+ [-+]



Systems biology?



! Rapamycin has no 100% activity  $\rightarrow$  S6K1  $\downarrow$   $\rightarrow$  PI3K  $\uparrow$



$$\frac{d[K_1]}{dt} = -k_1[K_1][E_1] + k_2[K_1^*][E_2]$$

approximation

Measurements of concentrations  $\rightarrow$  establishment of model.  
 in the previous feedback loop system, inhibiting IRS, upstream of the cascade, gives a better control than rapamycin. And a combination is predicted to be even more efficient.

- specificity (target)
- therapeutic window

Gene Sequencing:

1995	<i>H. influenza</i> (bacteria)	1.8M	
1997	<i>S. cerevisiae</i> (yeast)	12M	$\leftarrow$ microarray
1998	<i>C. elegans</i>	100M	
1999	<i>D. melanogaster</i>	16.5M	
2000	Human	3000M	

Environmental Stress Response:  $\sim$  1000 out of 6,000 genes in cerevisiae change expression regardless of the external stimulus

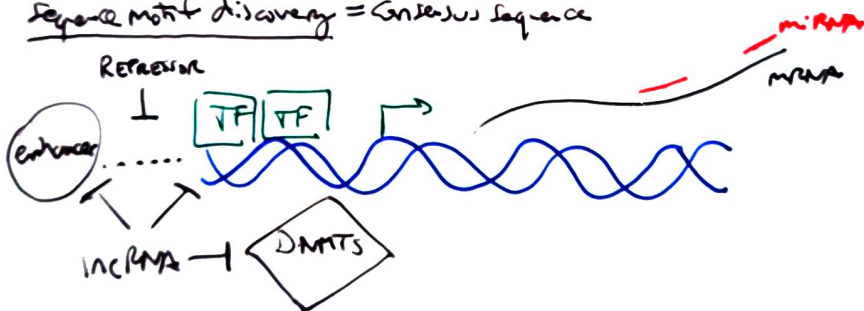
Gene Signature = set of genes useful for predicting something; e.g., disease.

$\downarrow$   
 gene ontology database (GO): Knowledge on function, location, etc.

! 90% of dataset not involved when perturbing a particular pathway

$\rightarrow$  transcription factors common to many gene expression

$\uparrow$   
Sequence motif discovery = consensus sequence



easy to measure!

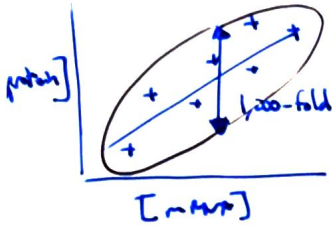
1999: reverse genetics (yeast) kindly/perturbation  
 = modify genome  $\rightarrow$  observe phenotype, systematically.

2000: CHIP-chip Chromatin immunoprecipitation  
 (2009) ↳ micro array protein-DNA interaction

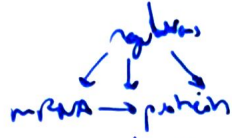
Muradap: CHIP-seq

2001-2002: protein-protein interaction / arrays

→ mass spectrometry



how much noise?



regulate true from a statistical level, but not very predictive

- reverse genetics, deletion/mutation, of one gene at a time  
 → identify a gene that is a master regulator for the regulation of gene expression in response to a stimulus

- RNA-seq → identify set of genes / pathways dysregulated upon stimulus

