



National Institutes of Health -- Center for Information Technology



Division of Computational Bioscience

# Lognormal Pattern

*of* Exon size distributions  
*in* Eukaryotic genomes

*Yaroslav Ryabov*

# Outline

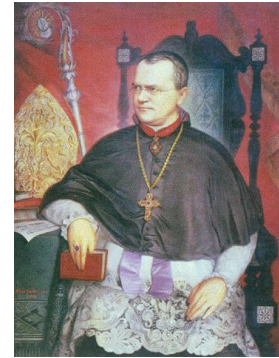
- Global and local approaches in analysis of genetic information
- Vocabulary of contemporary genetics:  
*Exons and Introns*
- What we can learn from exon size distributions?  
*Lognormal pattern and two classes of exons*
- How we can model exon size distributions of real genomes?
- What could be the biological reason for observed pattern of exons size distributions?
- Conclusions

# Analyzing Genetic Information

## Global approach

analyzing properties  
of entire genome

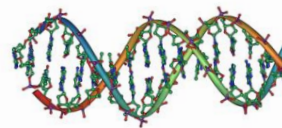
Recombination of  
inherited properties  
Frequency of mutations  
Size of genome *etc.*



*Georg Johann Mendel*  
1866

## Local approach

analyzing details of  
nucleotide sequence



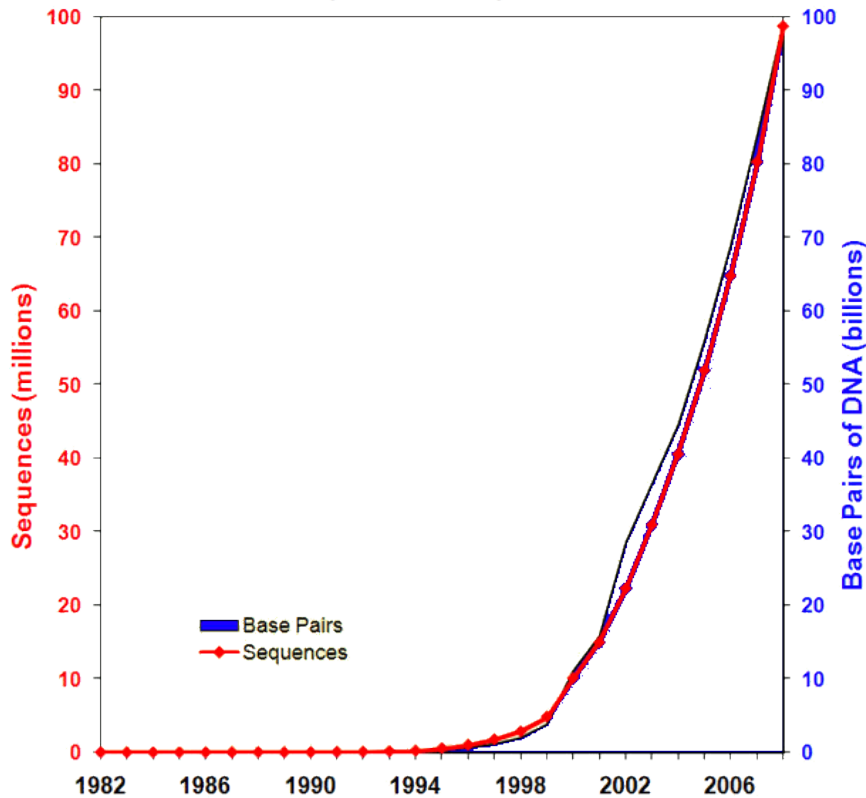
*Watson and Crick*  
1953  
*Marshall Nirenberg*  
1968



**Basic Local  
Alignment Search Tool**  
1990

# Analyzing Genetic Information in Post-Genomic Era

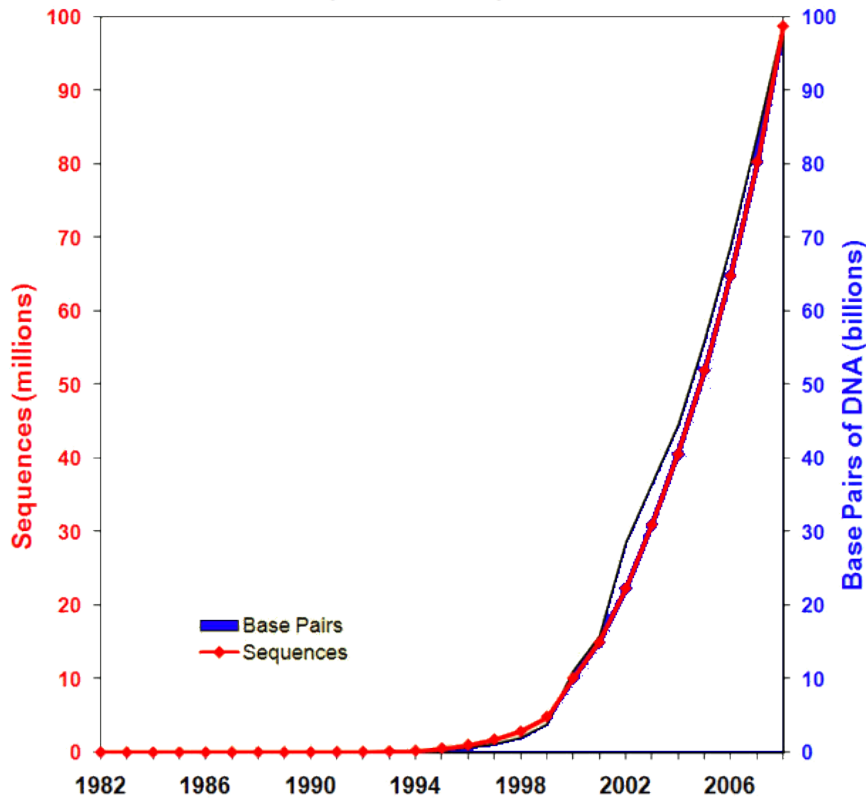
Growth of GenBank  
(1982 - 2008)



More than 60 animal  
genomes with complete  
annotations

# Analyzing Genetic Information in Post-Genomic Era

Growth of GenBank  
(1982 - 2008)



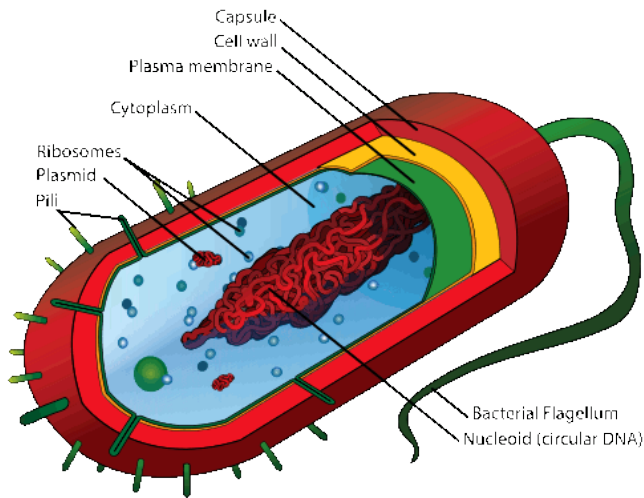
More than 60 animal  
genomes with complete  
annotations

**Back to Global Approach ?**

# Prokaryote

pro + karyon

*before + nucleos*

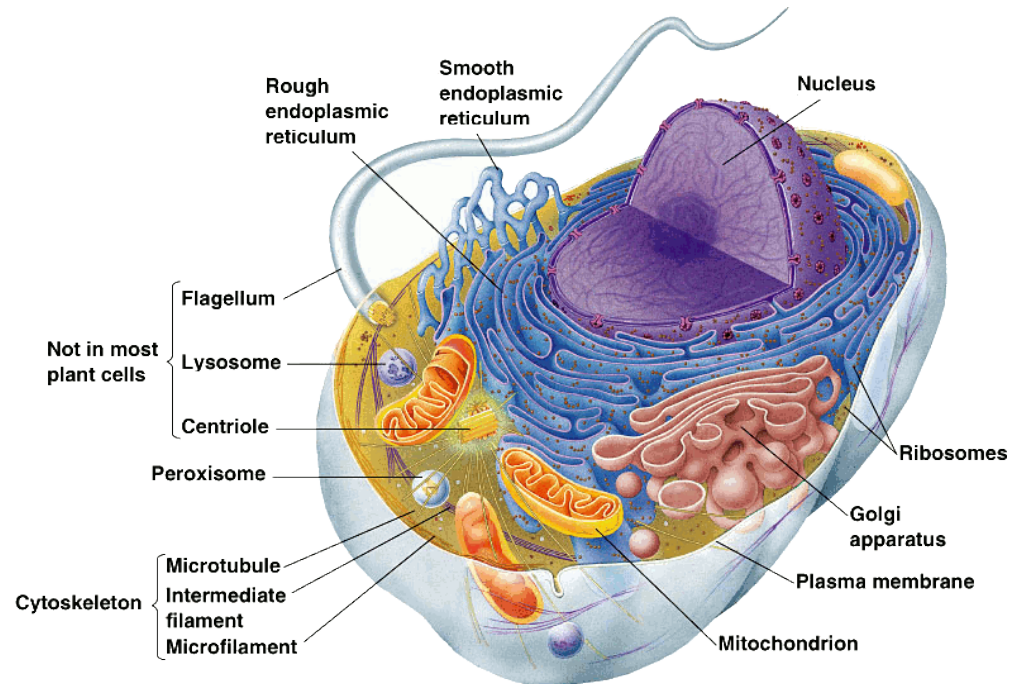


and

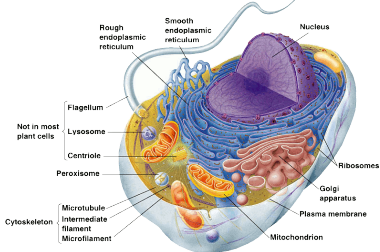
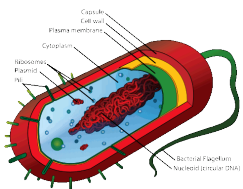
# Eukaryote

eu + karyon

*good + nucleos*

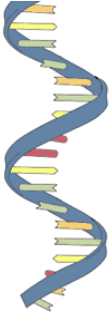
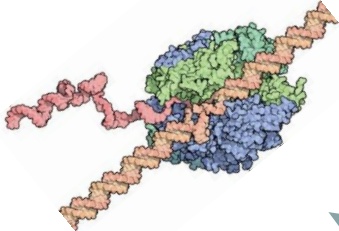


# Gene expression

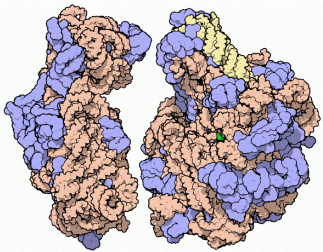


## mRNA

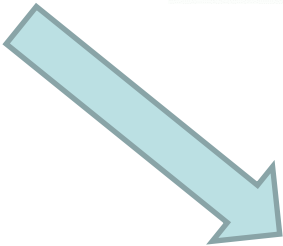
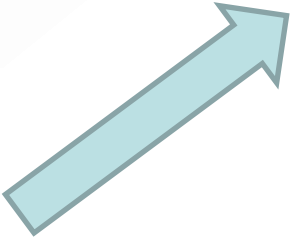
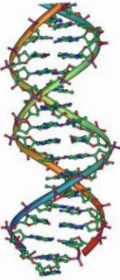
Transcription with RNA polymerase



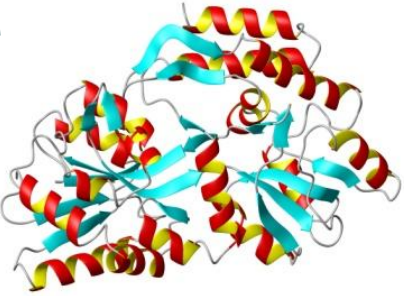
Translation with ribosome



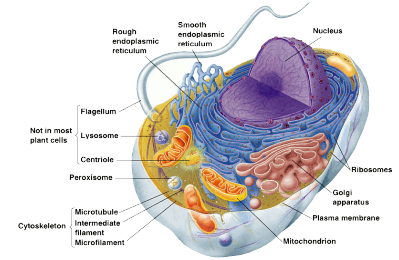
## DNA



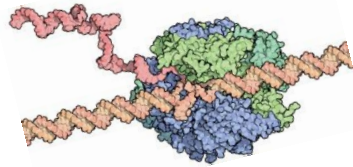
## Protein



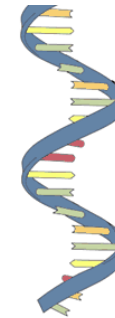
# Splicing in Eukaryotes



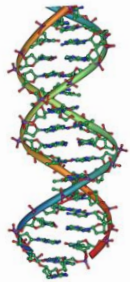
Transcription  
with RNA  
polymerase



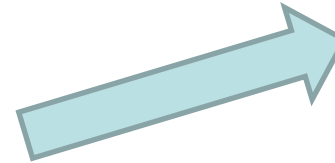
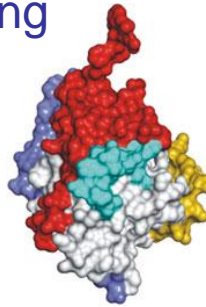
mRNA



DNA



Splicing



**Exons**



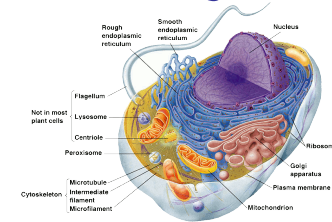
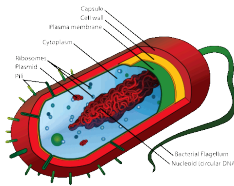
**Introns**



# Prokaryote

VS

# Eukaryote



Most of DNA code is used to produce some cellular product

Substantial fraction “silent” DNA regions



**Exons**

**Introns**

Short genomes: ~ 1 000 exons

Long genomes: ~ 100 000 exons

Long exons: ~ 1 000 base pairs

Short exons: ~ 100 base pairs

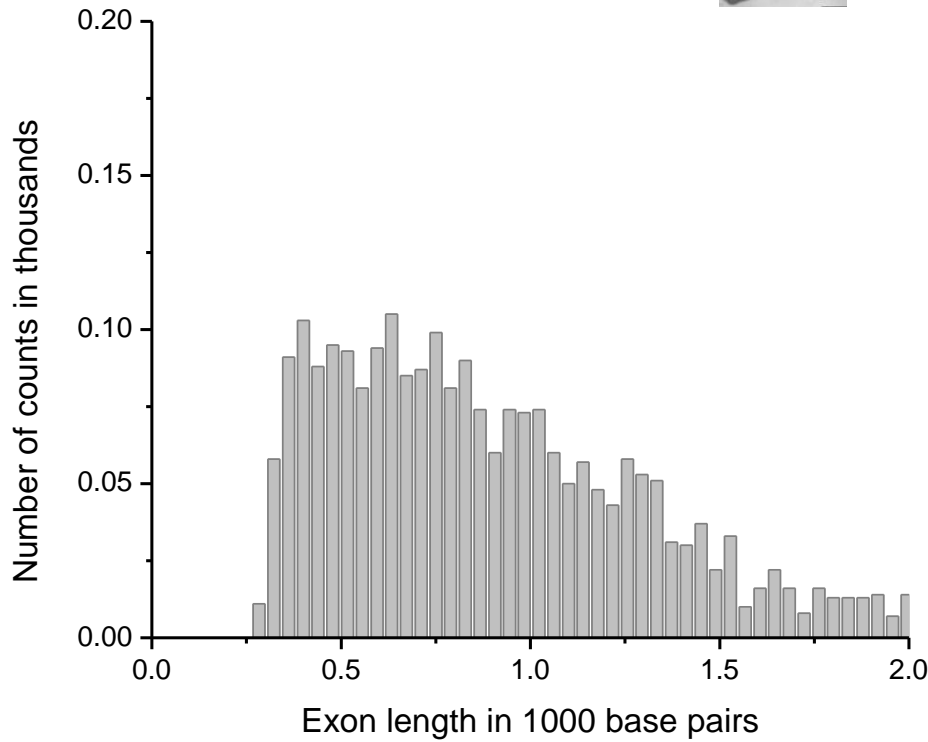
# Prokaryote

vs

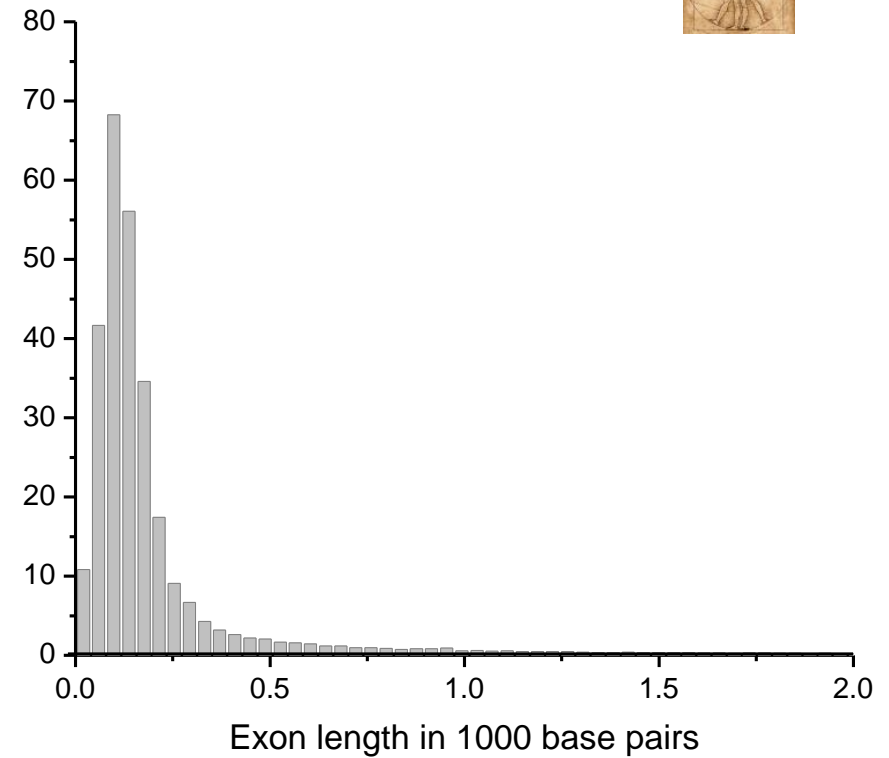
# Eukaryote

## Exon size distributions

Flavobacteria  
Bacterium

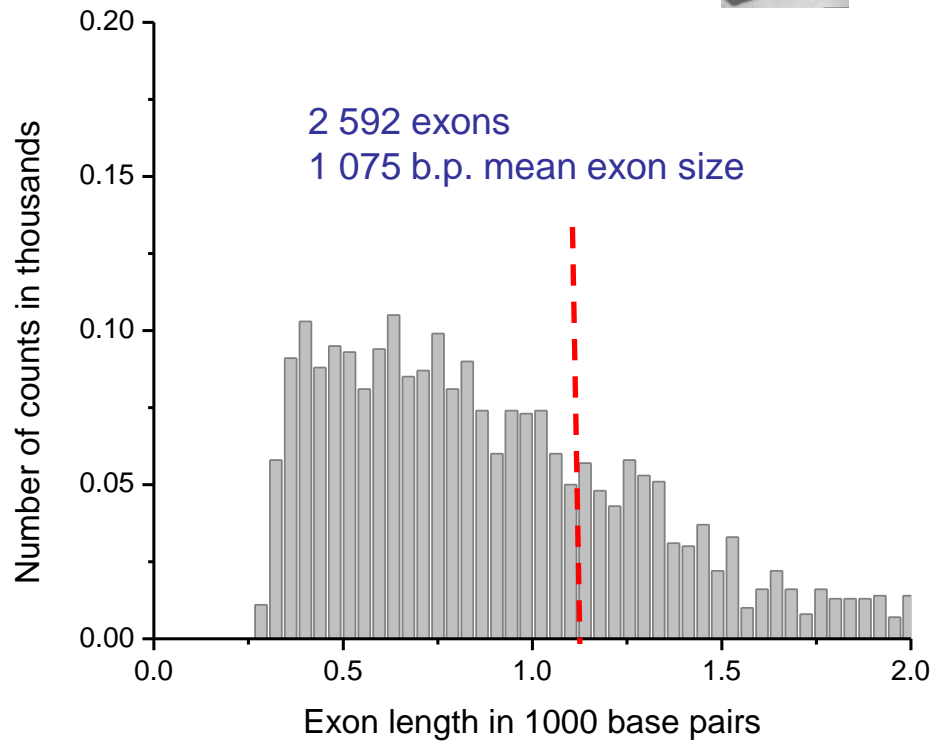


Homo Sapiens

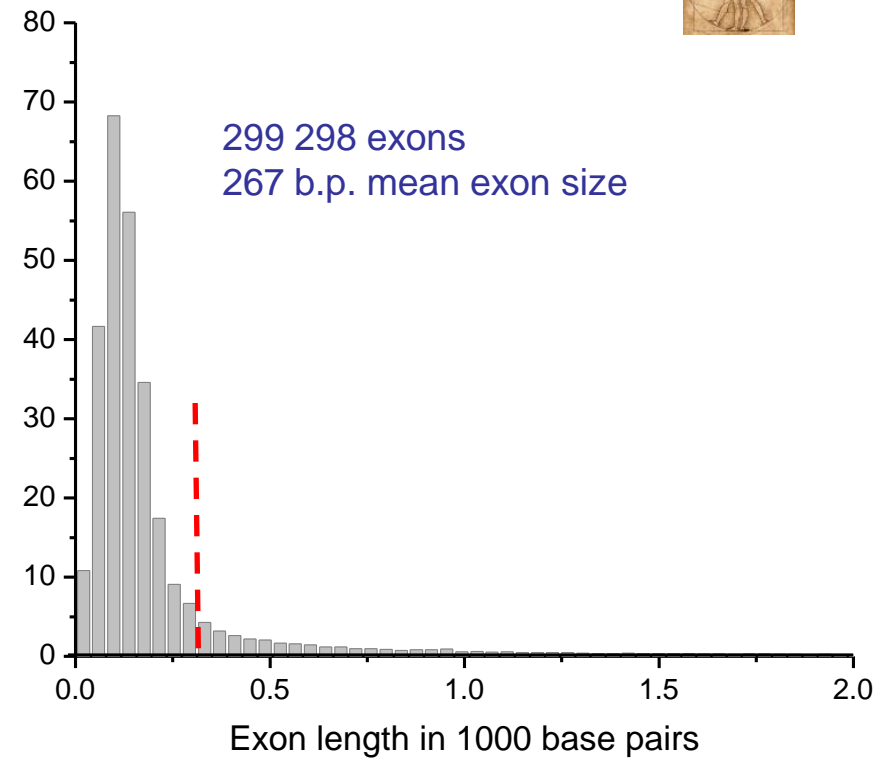


# What we can learn from exon size distributions ?

Flavobacteria  
Bacterium

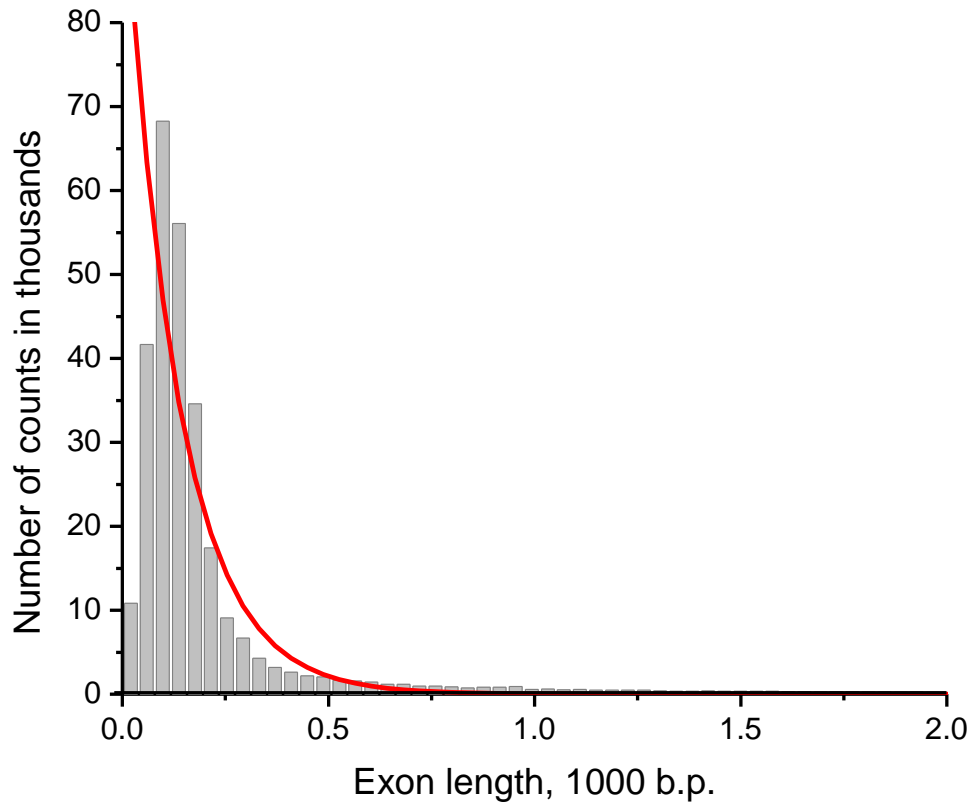


Homo Sapiens



# Poisson Process

$\lambda = \text{Const}$  frequency of splitting



## Assumption:

Probability to split DNA at any location and any time is *Constant*

## Consequence:

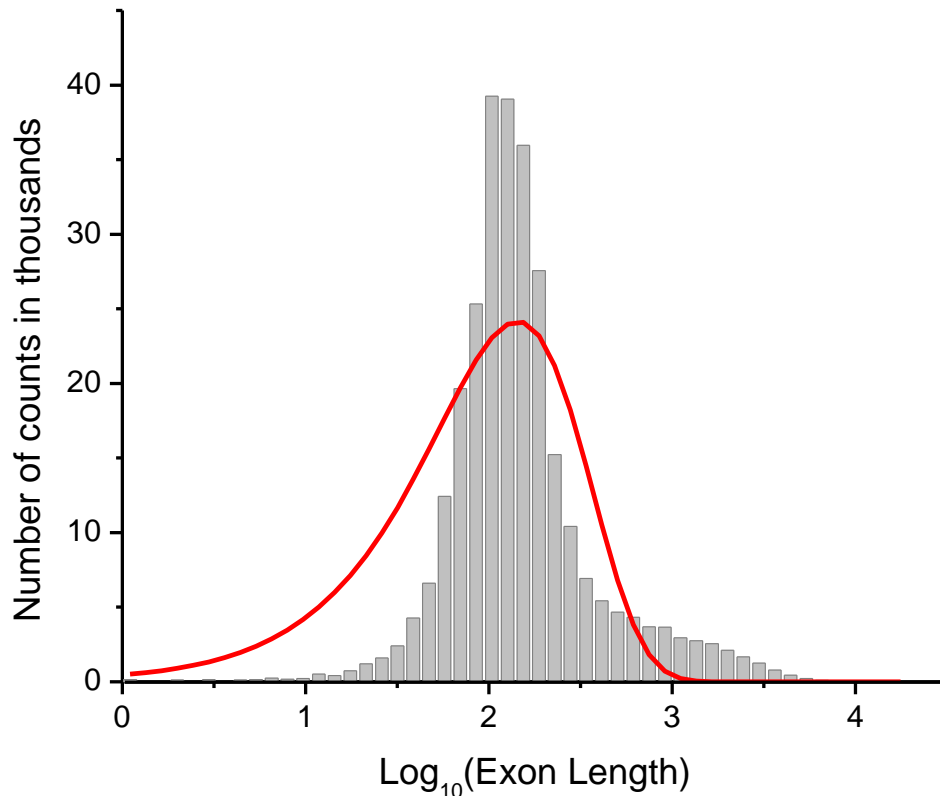
The lengths of intervals between splitting points obey

**Exponential** distribution

$$\frac{dN}{dE} = n_0 \lambda e^{-\lambda E}$$

# Poisson Process

$\lambda = \text{Const}$  frequency of splitting



## Assumption:

Probability to split DNA at any location and any time is *Constant*

## Consequence:

The lengths of intervals between splitting points obey

## Exponential distribution

$$\frac{dN}{dE} = n_0 \lambda e^{-\lambda E}$$

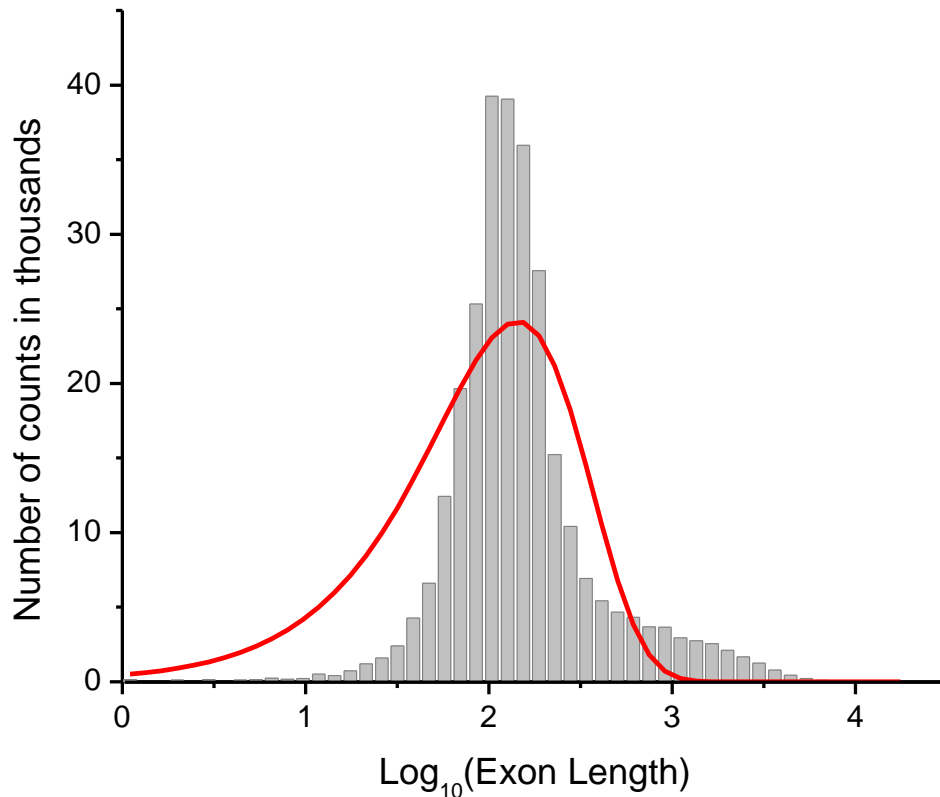
or in logarithm scale

$$\frac{dN}{d \log E} = n_0 \lambda e^{\log E - \lambda e^{\log E}}$$

# Poisson Process

Exponential distribution

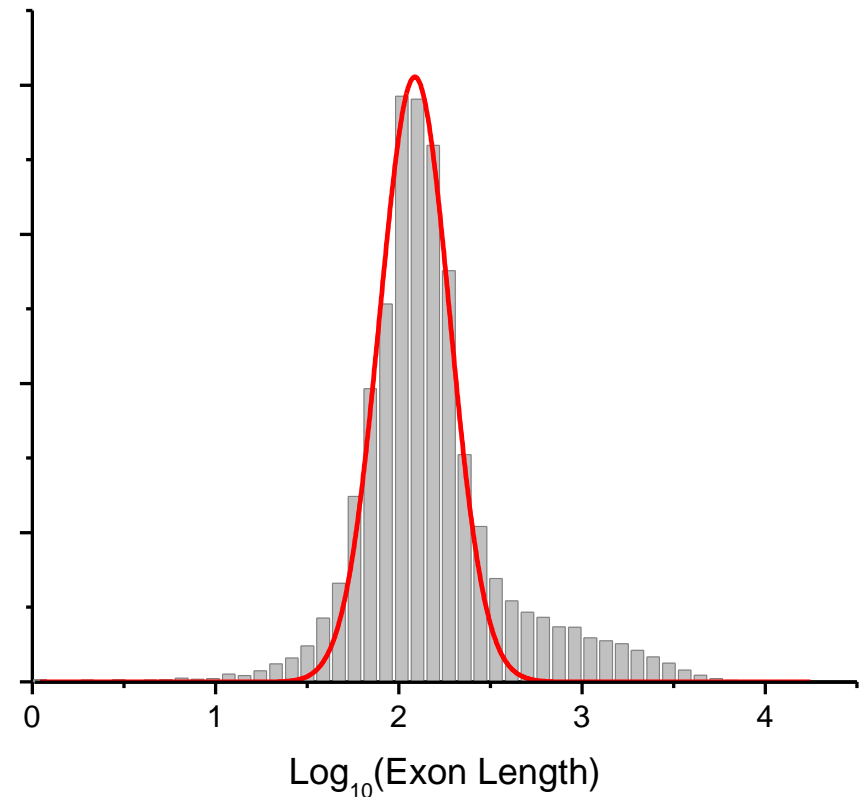
$$\frac{dN}{d \log E} = n_0 \lambda e^{\log E - \lambda e^{\log E}}$$



# Kolmogoroff Process

Lognormal distribution

$$\frac{dN}{d \log E} = \frac{n_0}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\log E - M)^2}{2\sigma^2}}$$



# Lognormal distribution

is a consequence of **Central Limit Theorem**

## Sum

of random variables

$$\Sigma = \xi_1 + \xi_2 + \xi_3 + \xi_4 + \dots$$

## Normal distribution

$$\frac{dP}{d\Sigma} = \frac{1}{\sqrt{2\pi\sigma_\Sigma^2}} e^{-\frac{(\Sigma - \langle \Sigma \rangle)^2}{2\sigma_\Sigma^2}}$$

## Product

of random variables

$$\Pi = \eta_1 \cdot \eta_2 \cdot \eta_3 \cdot \eta_4 \cdot \dots$$

## Lognormal distribution

$$\frac{dP}{d \log \Pi} = \frac{1}{\sqrt{2\pi\sigma_\Pi^2}} e^{-\frac{(\log \Pi - \langle \log \Pi \rangle)^2}{2\sigma_\Pi^2}}$$

$$\log \Pi = \log \eta_1 + \log \eta_2 + \log \eta_3 + \log \eta_4 + \dots$$

# Kolmogoroff process (1941)

$i = 0$



$E_0$



# Kolmogoroff process (1941)

$i = 0$



$E_0$

$i = 1$



# Kolmogoroff process (1941)



$E_0$  **Assumption:**  
Probability to split any exon is  
*Independent of exon size*

# Kolmogoroff process (1941)

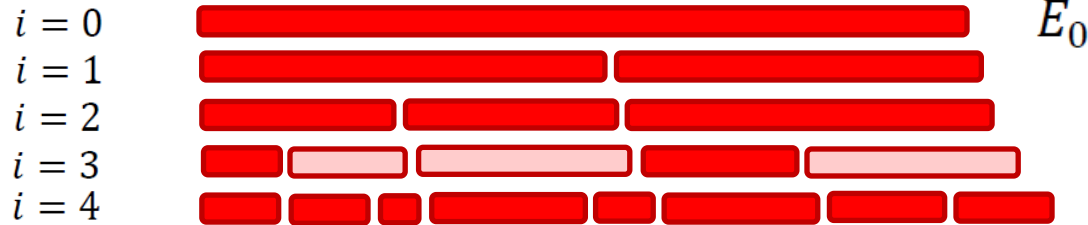


$E_0$

## Assumption:

Probability to split any exon is  
*Independent of exon size*

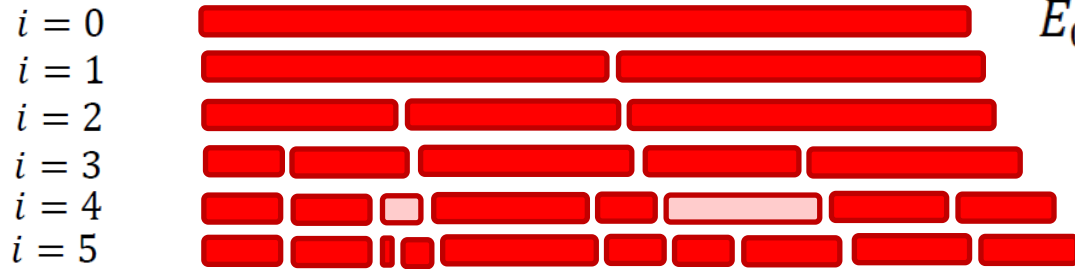
# Kolmogoroff process (1941)



## Assumption:

Probability to split any exon is  
*Independent of exon size*

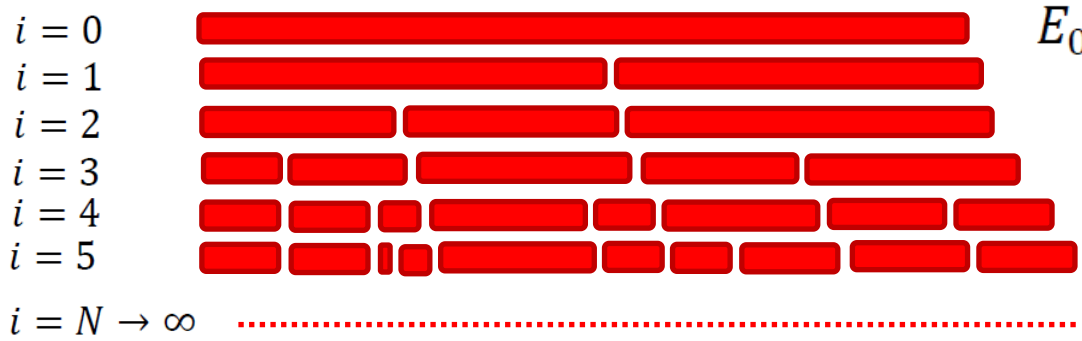
# Kolmogoroff process (1941)



## Assumption:

Probability to split any exon is  
*Independent of exon size*

# Kolmogoroff process (1941)



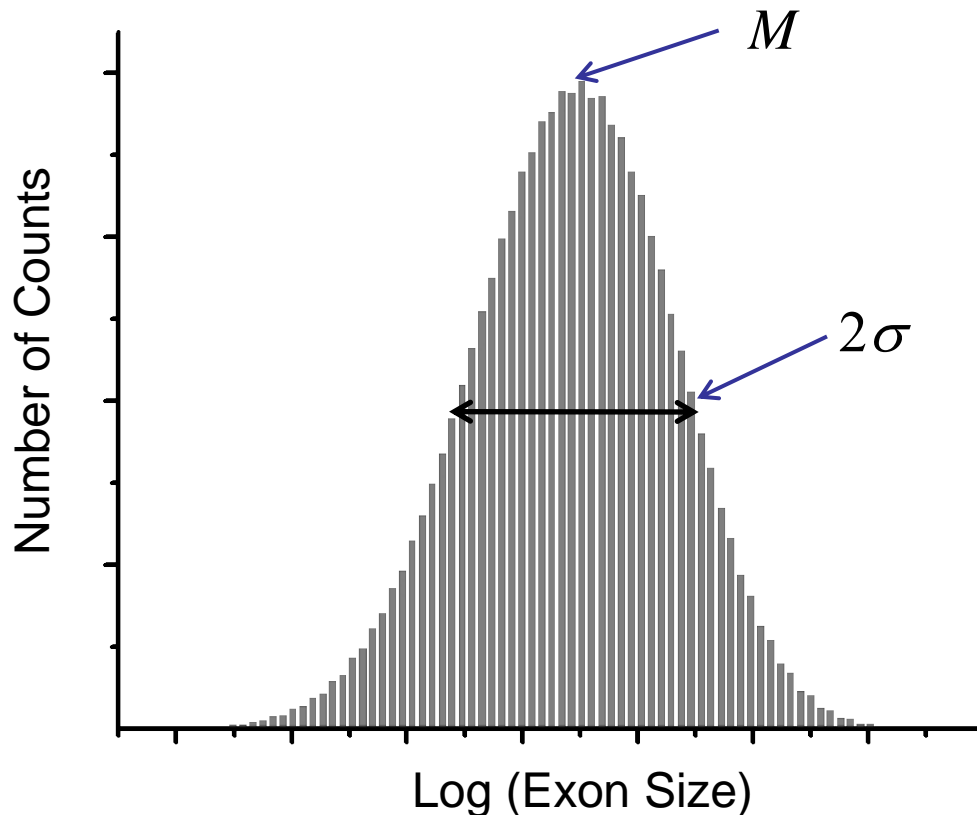
## Assumption:

Probability to split any exon is  
*Independent of exon size*

## Consequence:

The lengths of exons obey

**Lognormal** distribution

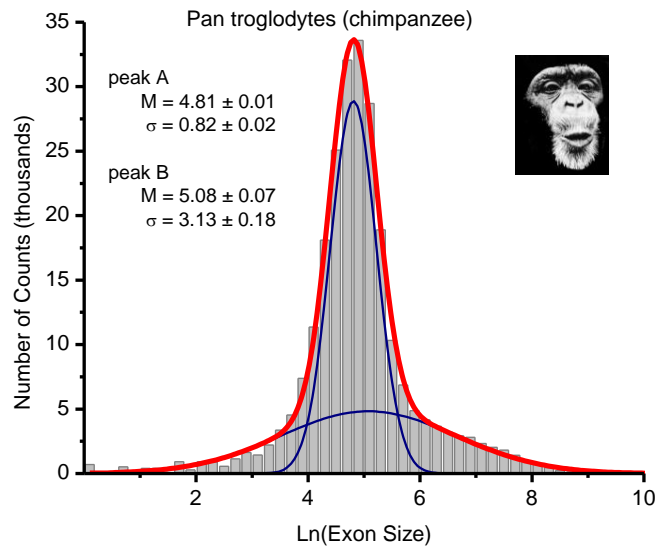
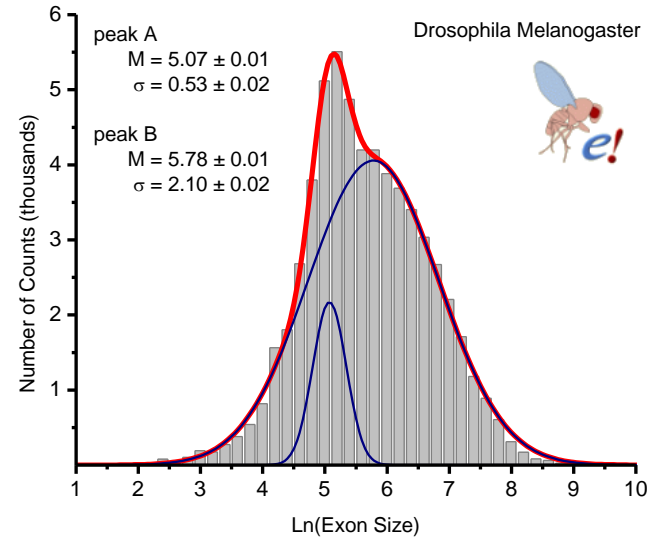
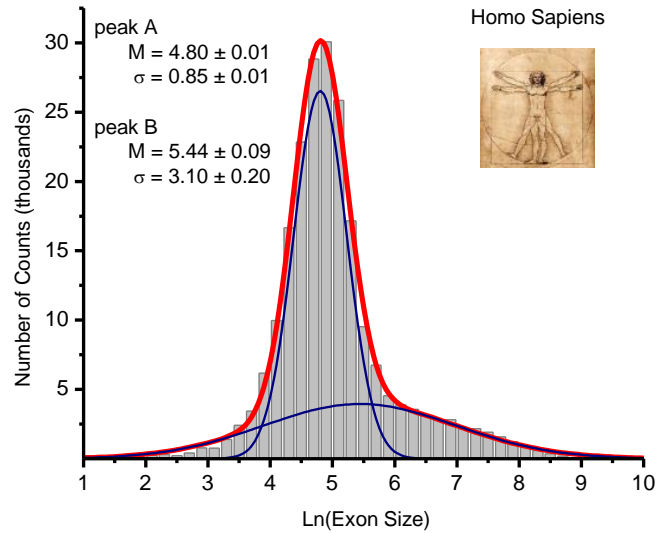


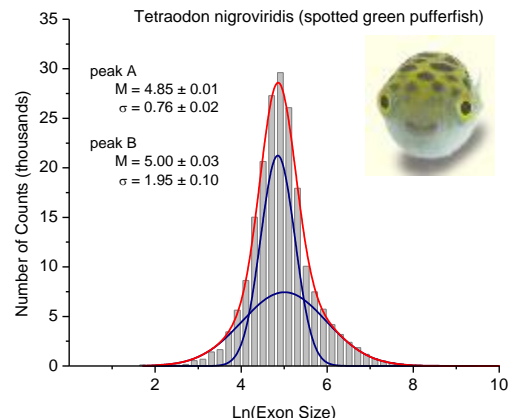
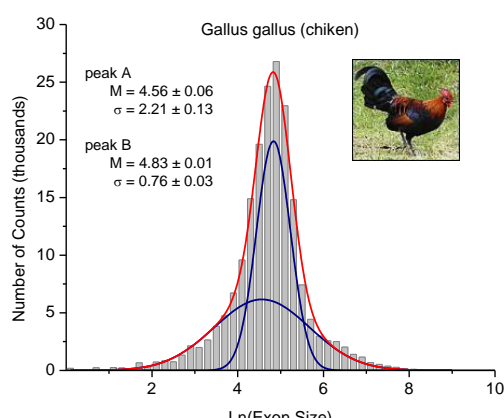
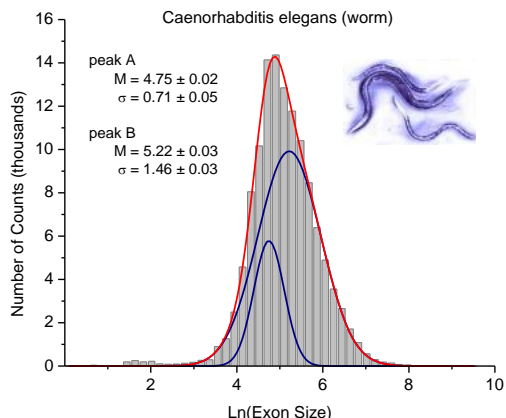
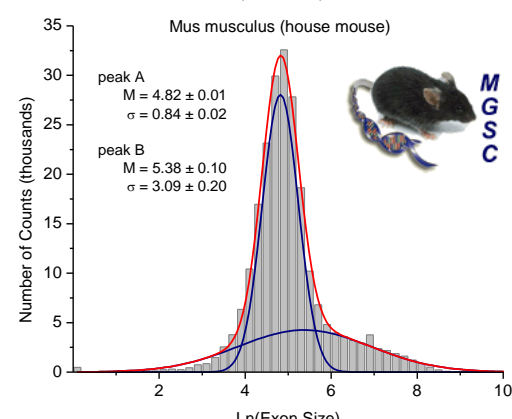
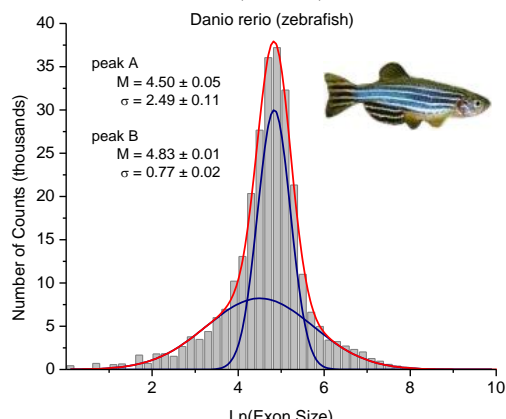
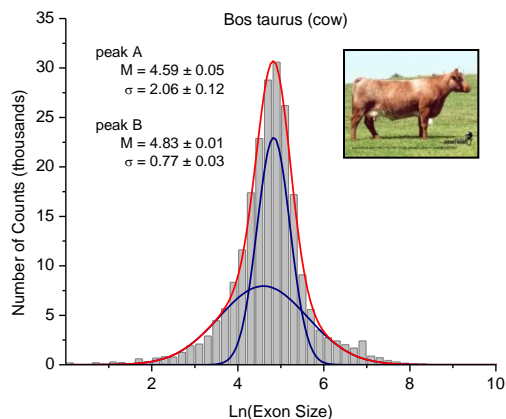
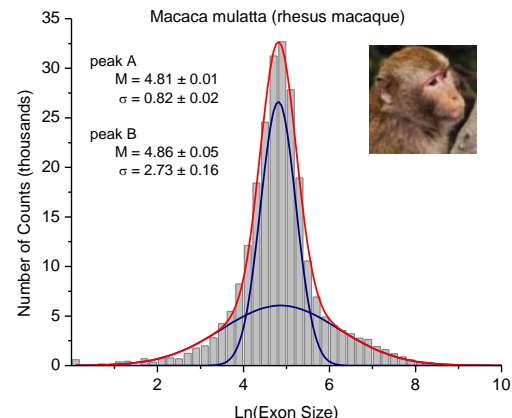
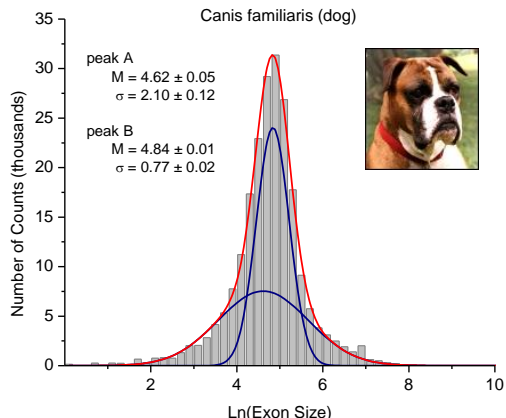
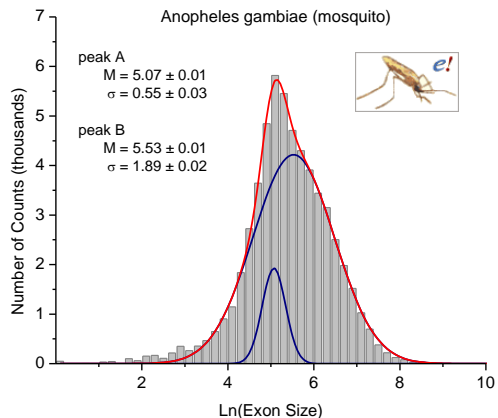
$$\frac{dP}{d \log E} = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\log E - M)^2}{2\sigma^2}}$$

$$\log E_0 - M \sim N$$

$$\sigma^2 \sim N$$

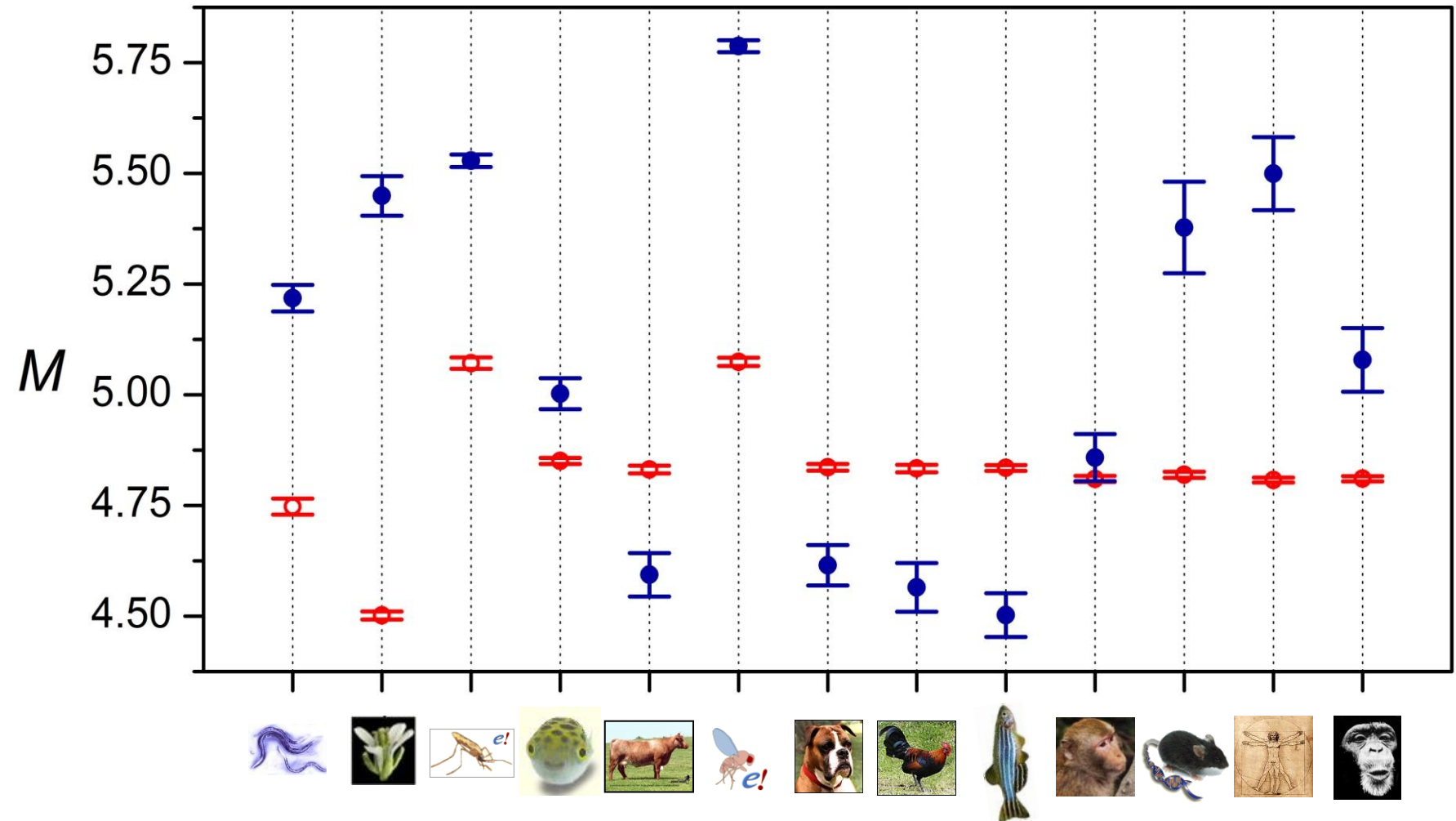
# Real Genomes: Two Lognormal Peaks



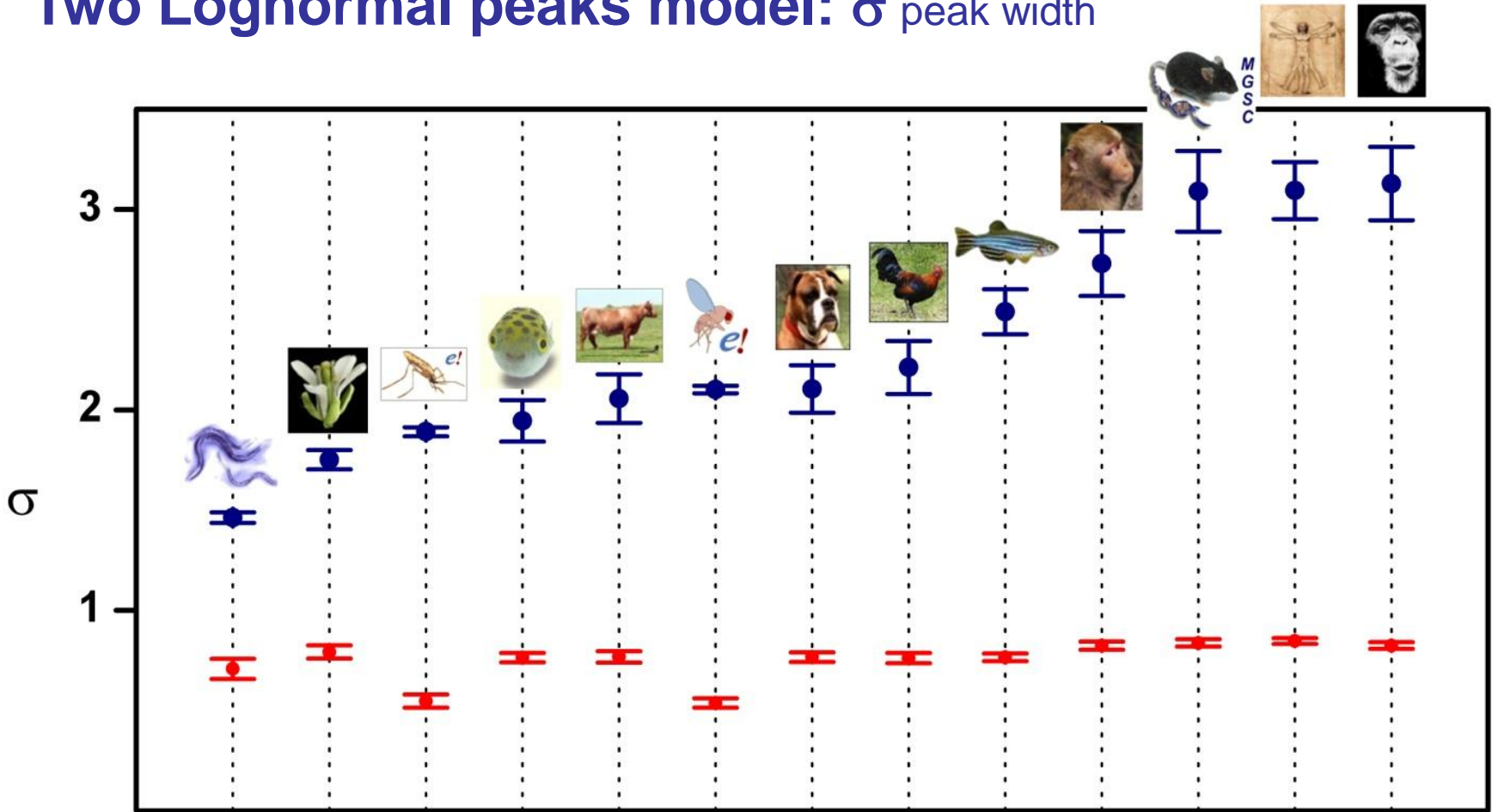




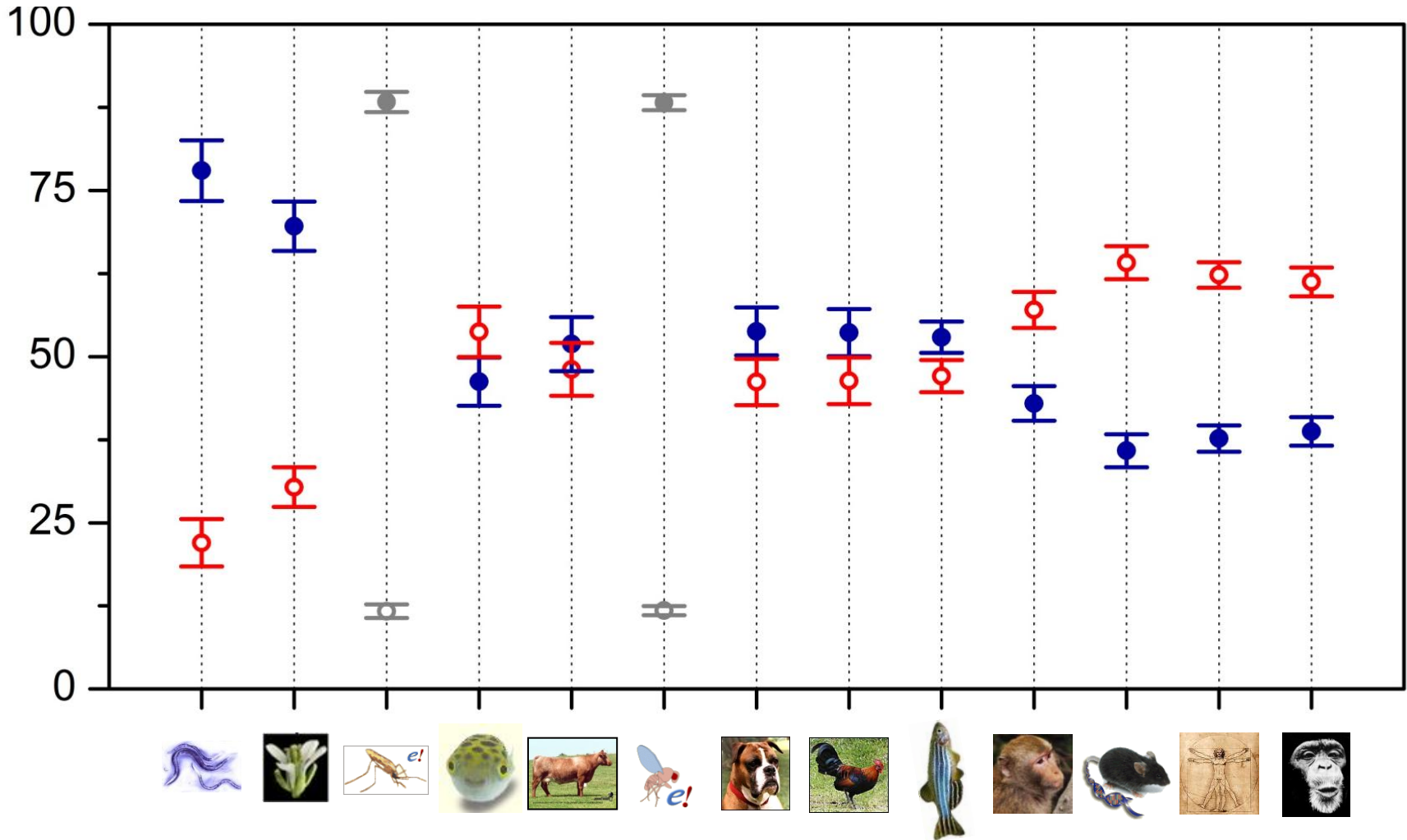
# Parameters of Two Lognormal peaks model: $M$ location of peak maximum



# Parameters of Two Lognormal peaks model: $\sigma$ peak width



# Parameters of Two Lognormal peaks model: peak area **Narrow** and Wide, %



## Summary of Observed Facts

- Exon lengths distributions of studied eukaryotic genomes can be fitted by Two Lognormal Distributions
- Parameters of those two peaks follow two distinctive patterns: changes of peak width and relative peak area correlate with complexity of species.
- This may indicate presence of two different classes of exons with different evolutionary histories

# **How we can model exon size distributions of real genomes ?**

**Growth of total genome length**

**Duplications in genome code**

**Merging exons together (intron loss)**

**Exon loss**

# Parameters of Kolmogoroff Process

as model parameters for elementary “exon modification” event

**If** any exon of length  $l_i$  has probability  $p$  to be modified during time interval  $\Delta t$

**and**

$Q(k)$  average number of exons with sizes  $l_{i+1} \leq k l_i$   $\int_0^\infty |\ln(k)|^3 dQ(k) < \infty$

**with**

$$A = \frac{1}{Q(\infty)} \int_0^\infty \ln(k) dQ(k) \quad B^2 = \frac{1}{Q(\infty)} \int_0^\infty (\ln(k) - A)^2 dQ(k)$$

**Then** the process converges to a lognormal distribution with

$$M(t) = M_0 + A \frac{p t}{\Delta t}$$

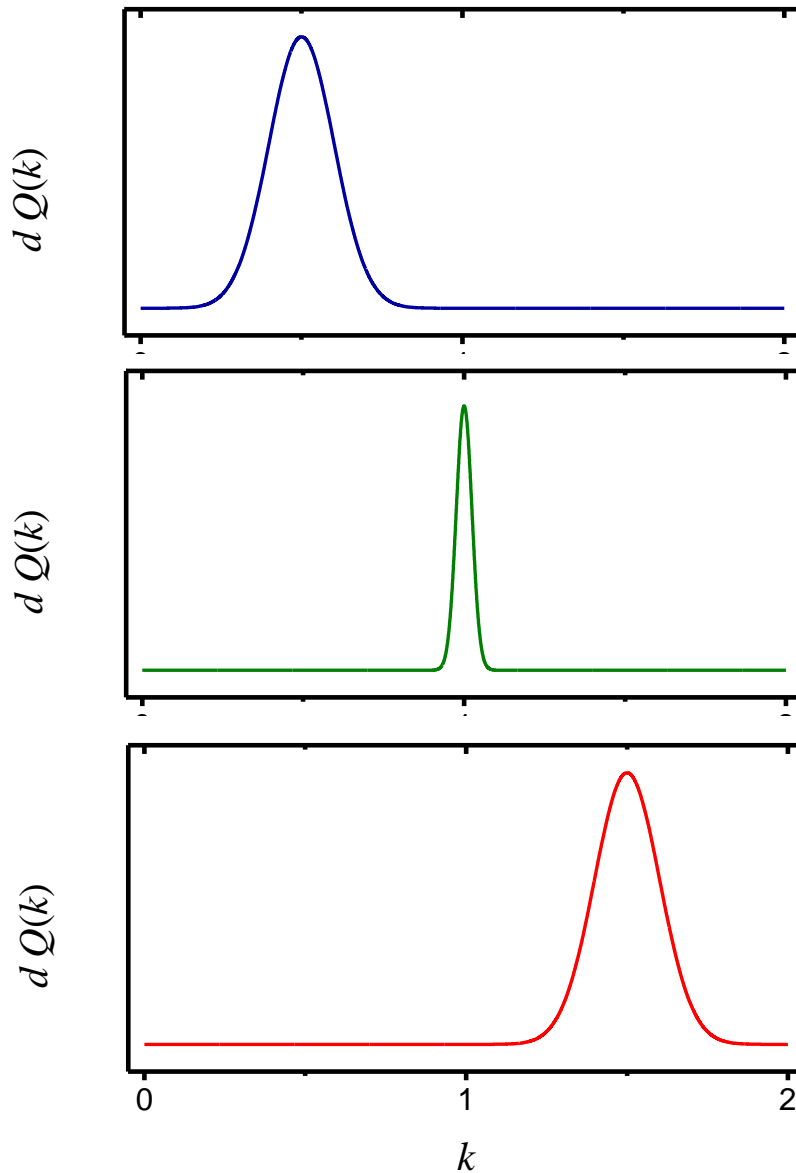
peak position

$$\sigma^2(t) = \sigma_0^2 + B^2 \frac{p t}{\Delta t}$$

peak width

# Parameters of Q(k) function

$dQ(k)$  distribution



$$A < 0$$

Exons splitting  
(decreasing Exon length)

$$A = \frac{1}{Q(\infty)} \int_0^{\infty} \ln(k) dQ(k)$$

$$A \approx 0$$

Exons duplicating

$$A > 0$$

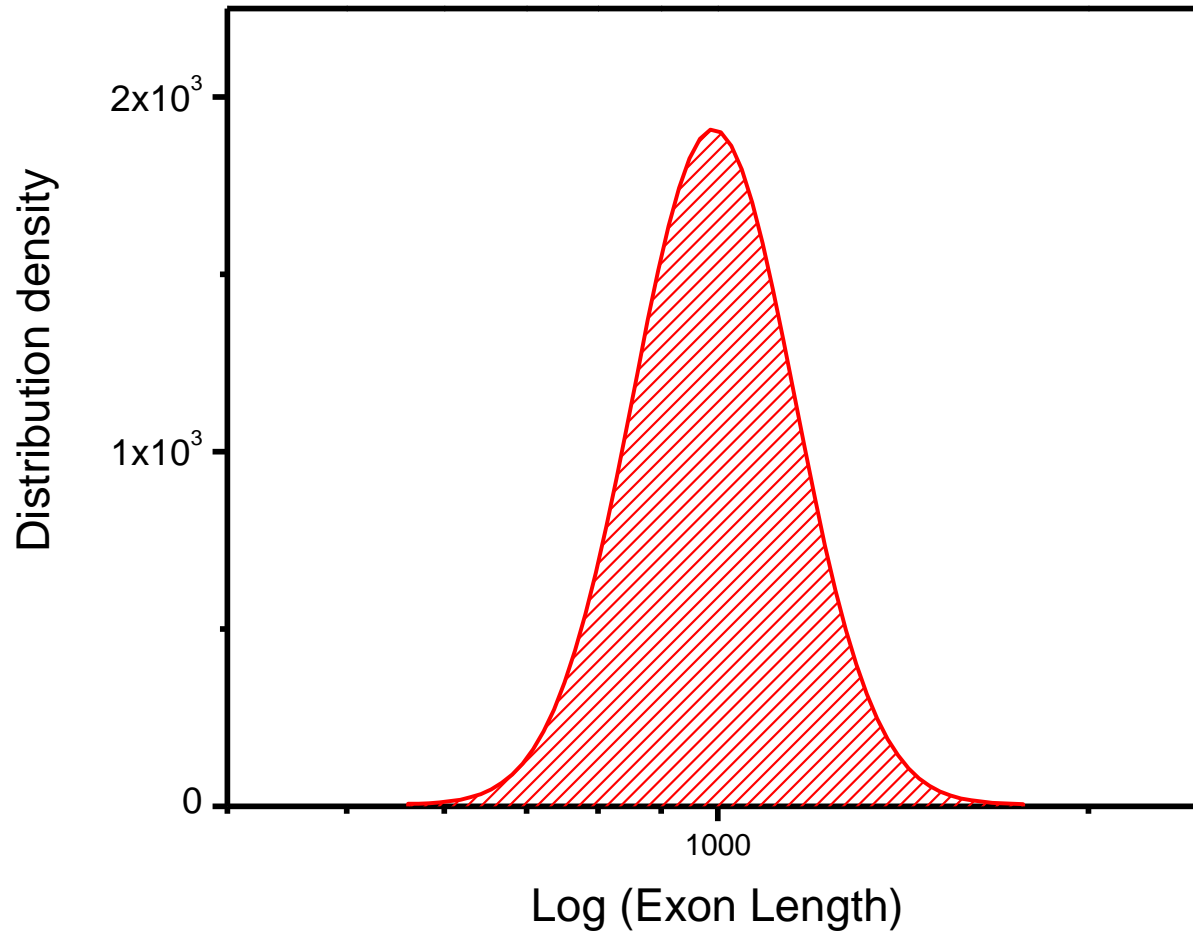
Increasing Exon length

# Modeling exon size distributions in real genomes

Initial exon size distribution

A mockup of a bacterial genome with 500 exons

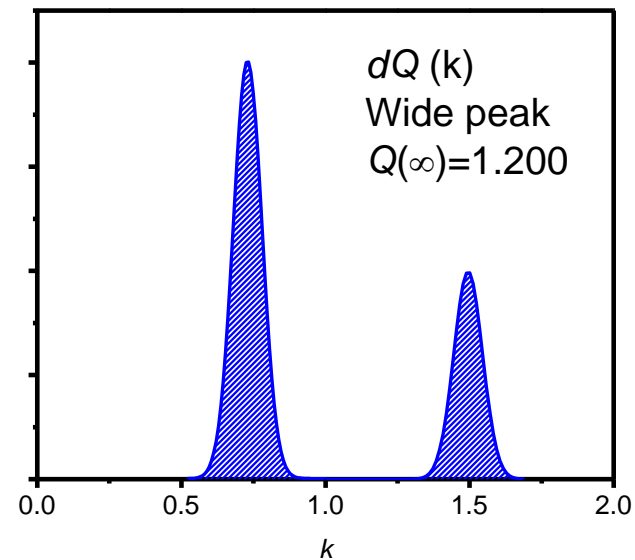
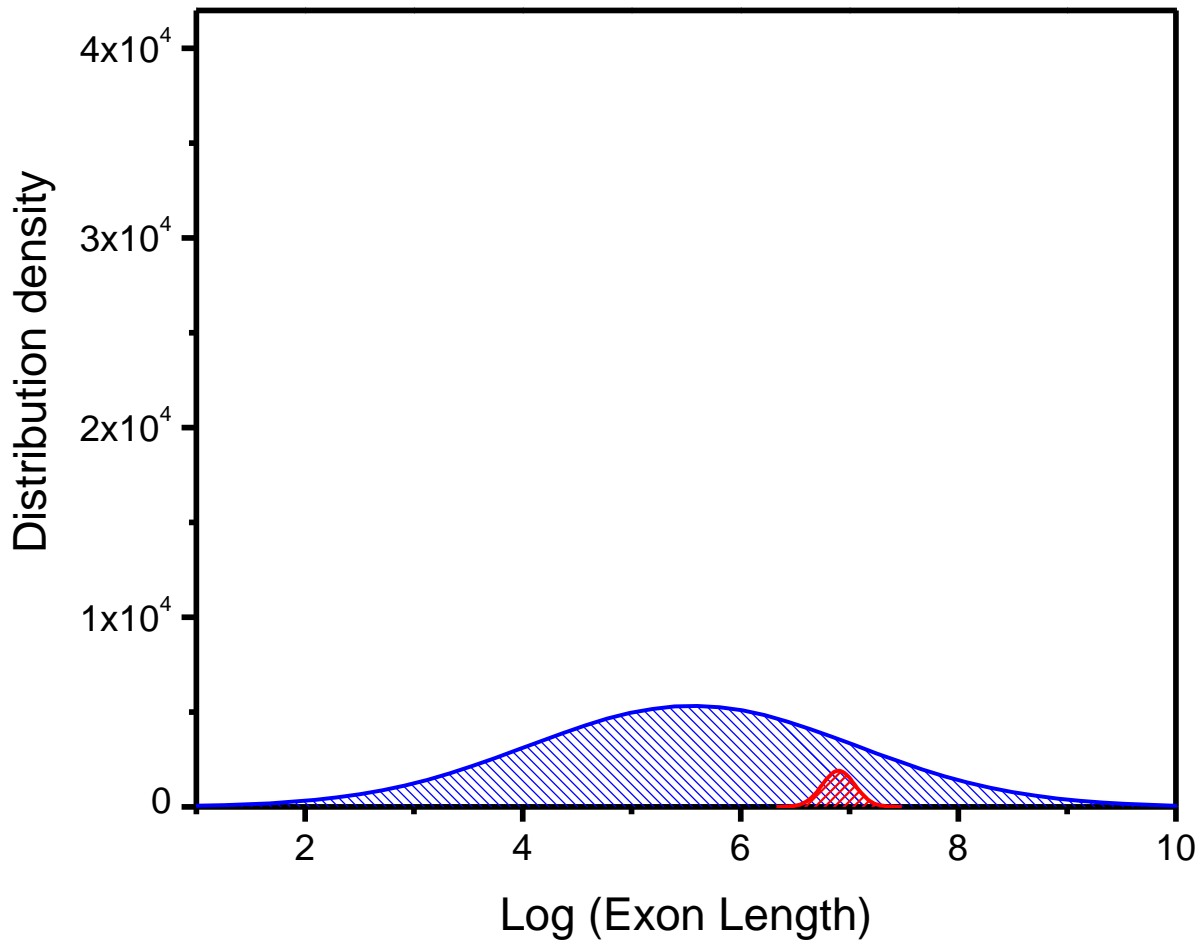
Having 1000 bp mean exons length





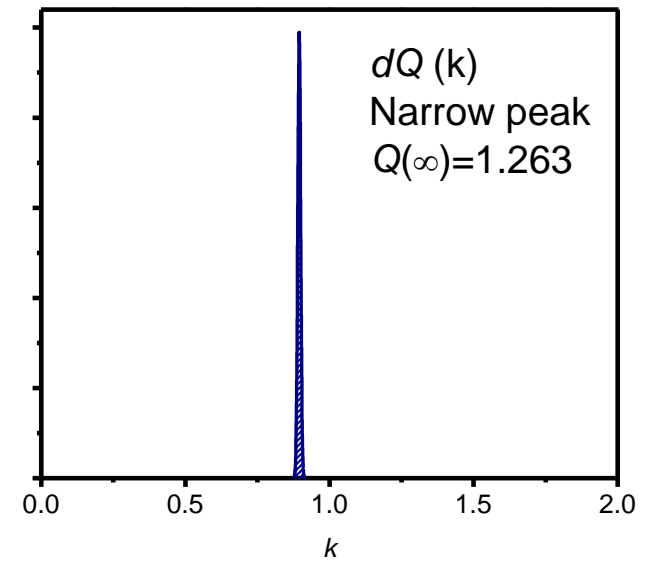
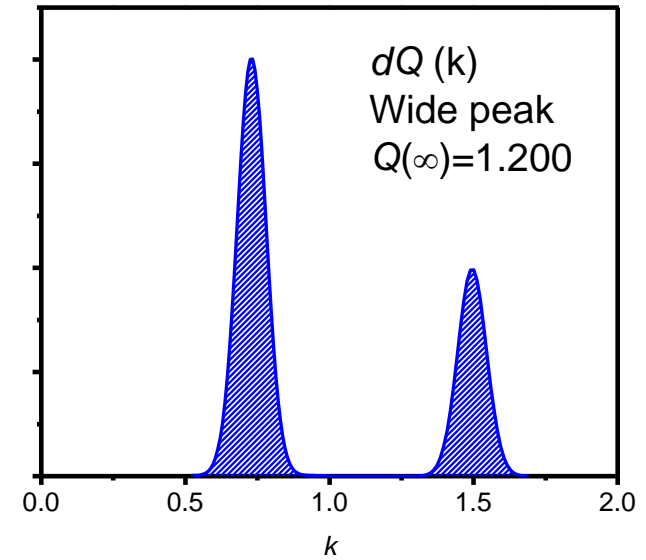
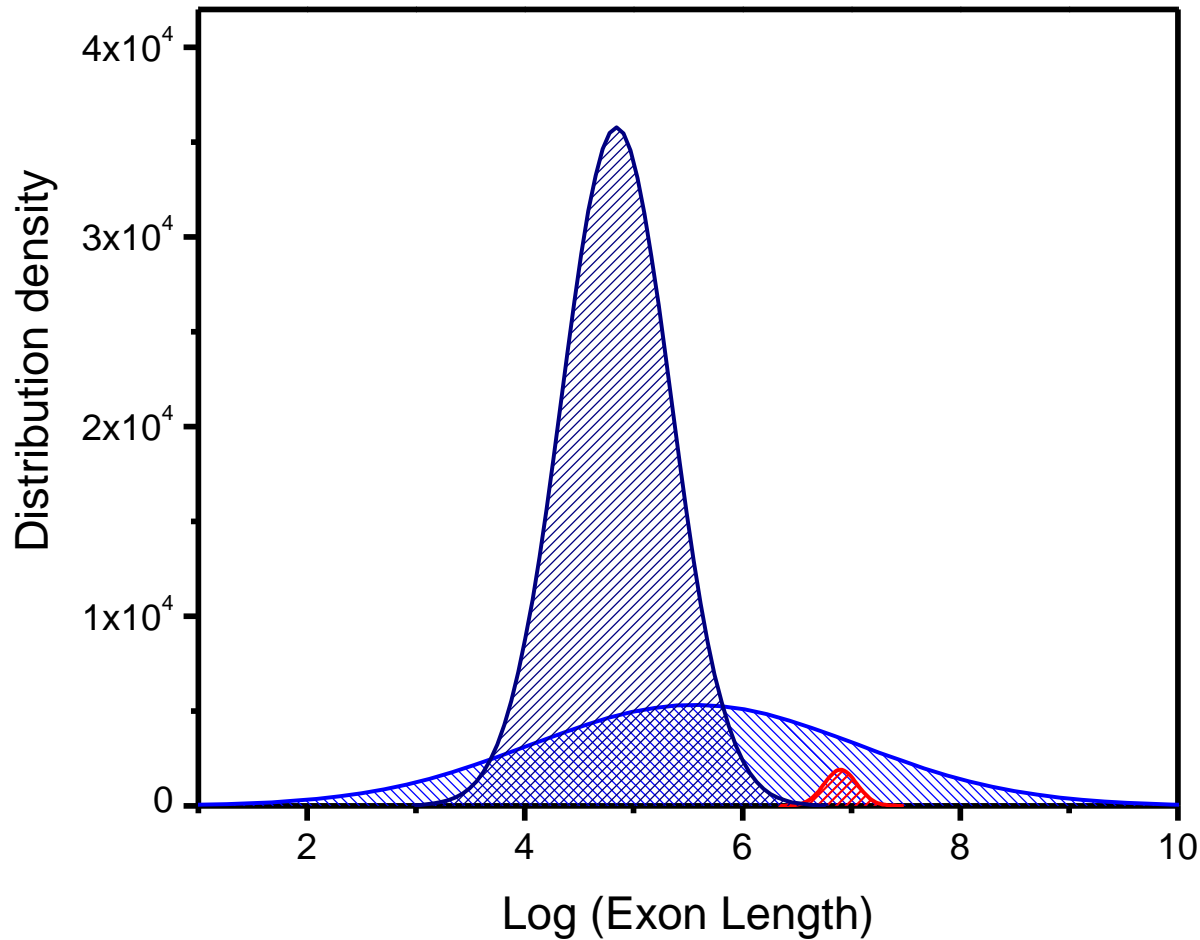
# Modeling exon size distributions in real genomes

15 000 time steps for  $p=0.001$



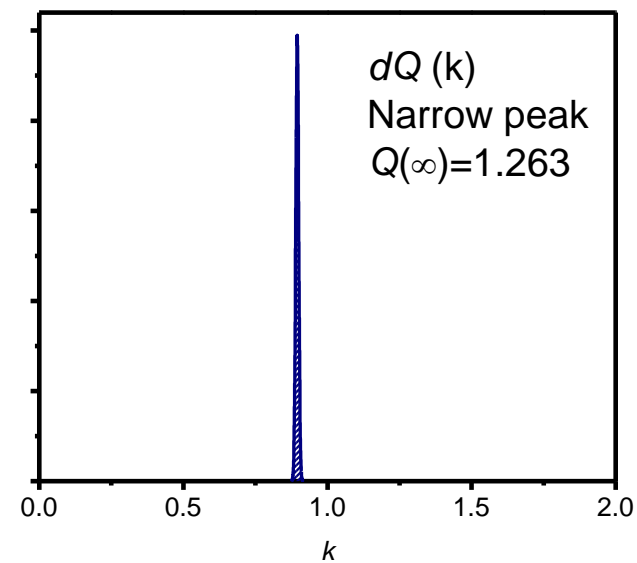
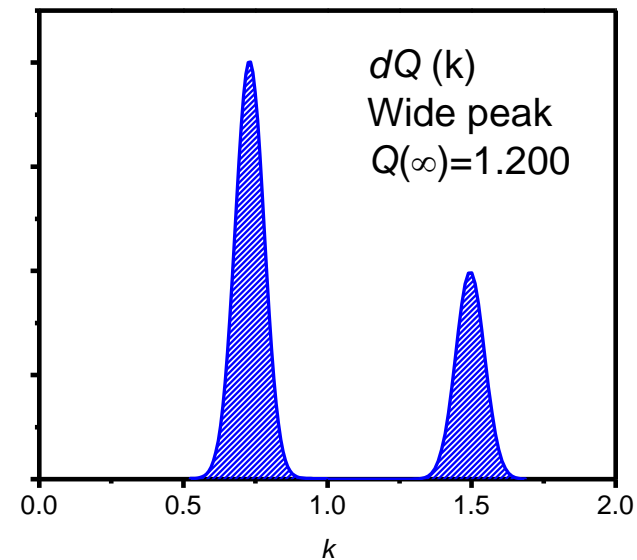
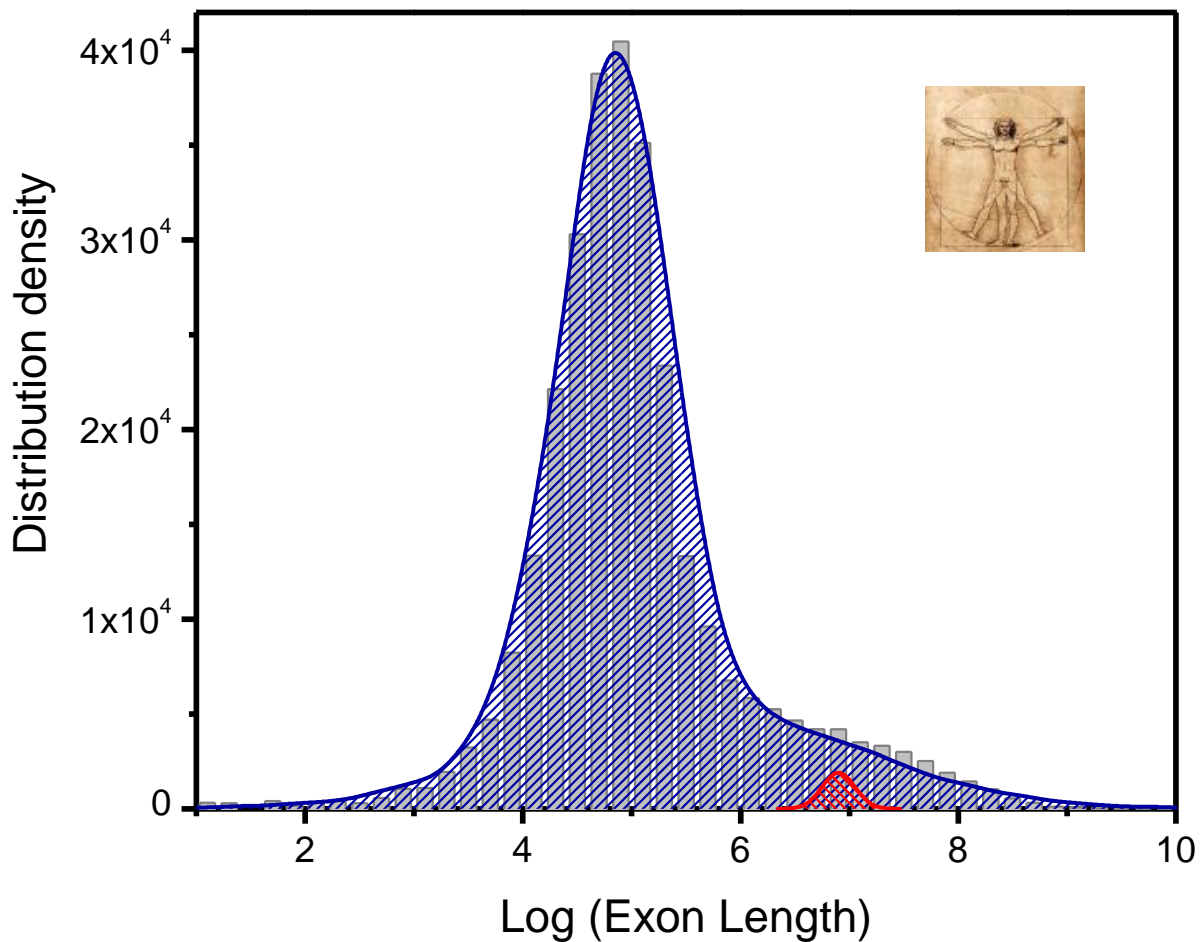
# Modeling exon size distributions in real genomes

15 000 time steps for  $p=0.001$



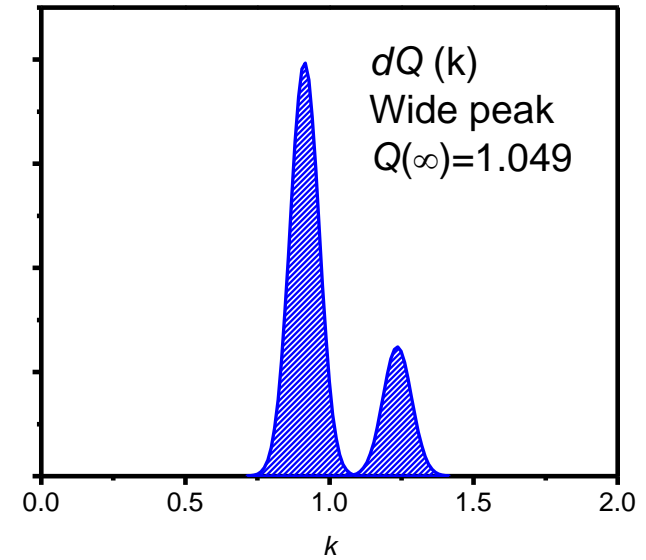
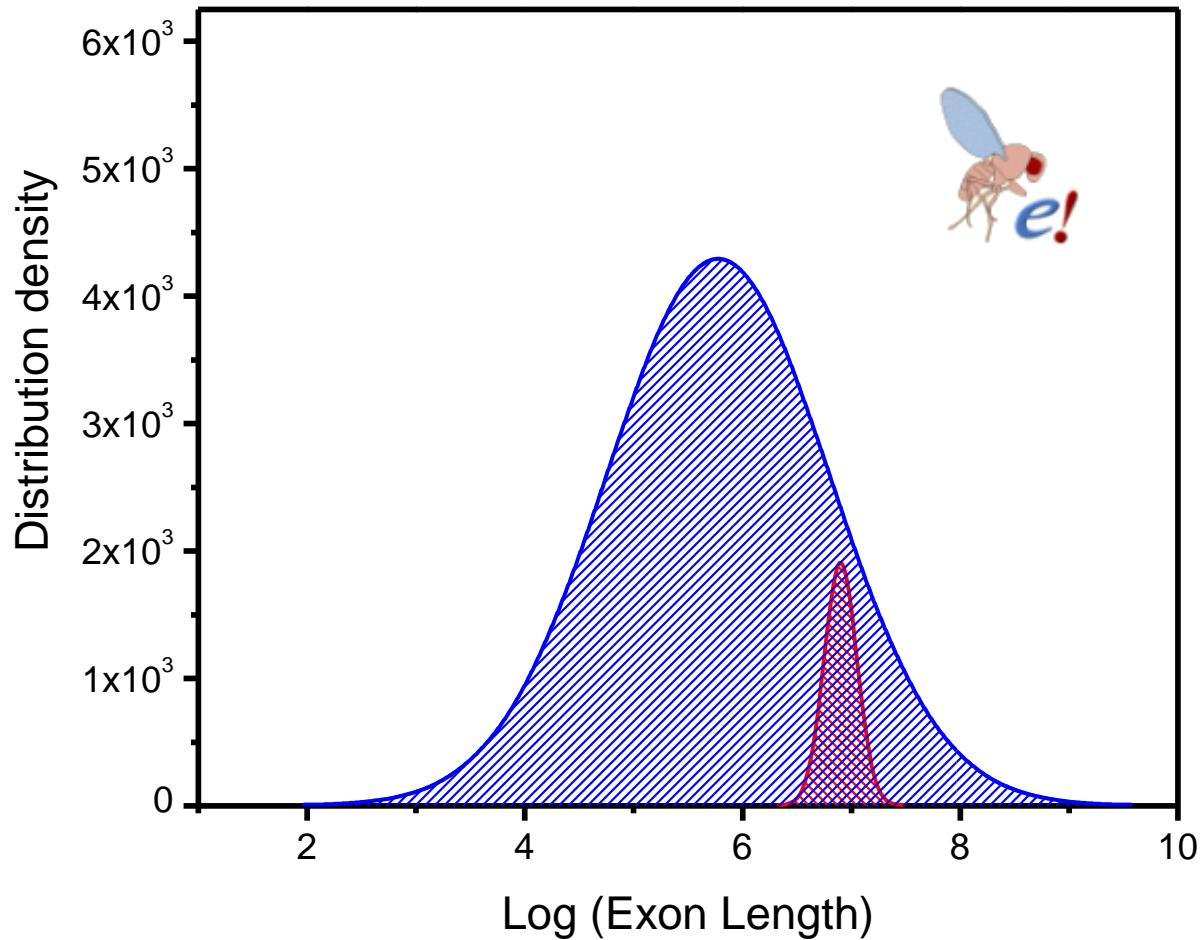
# Modeling exon size distributions in real genomes

15 000 time steps for  $p=0.001$



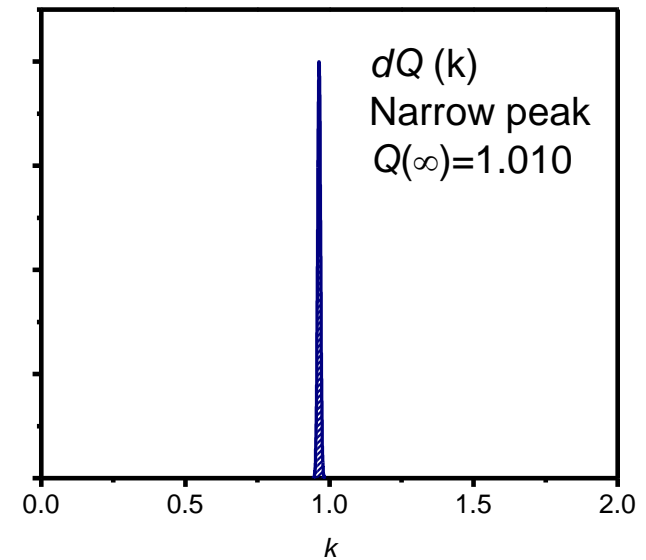
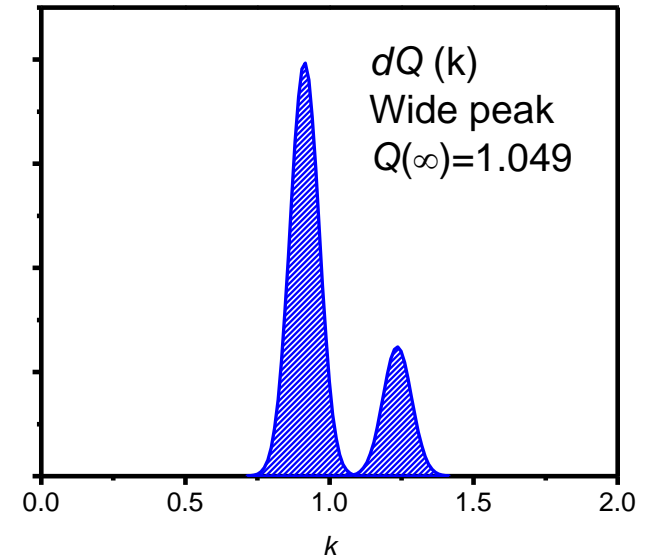
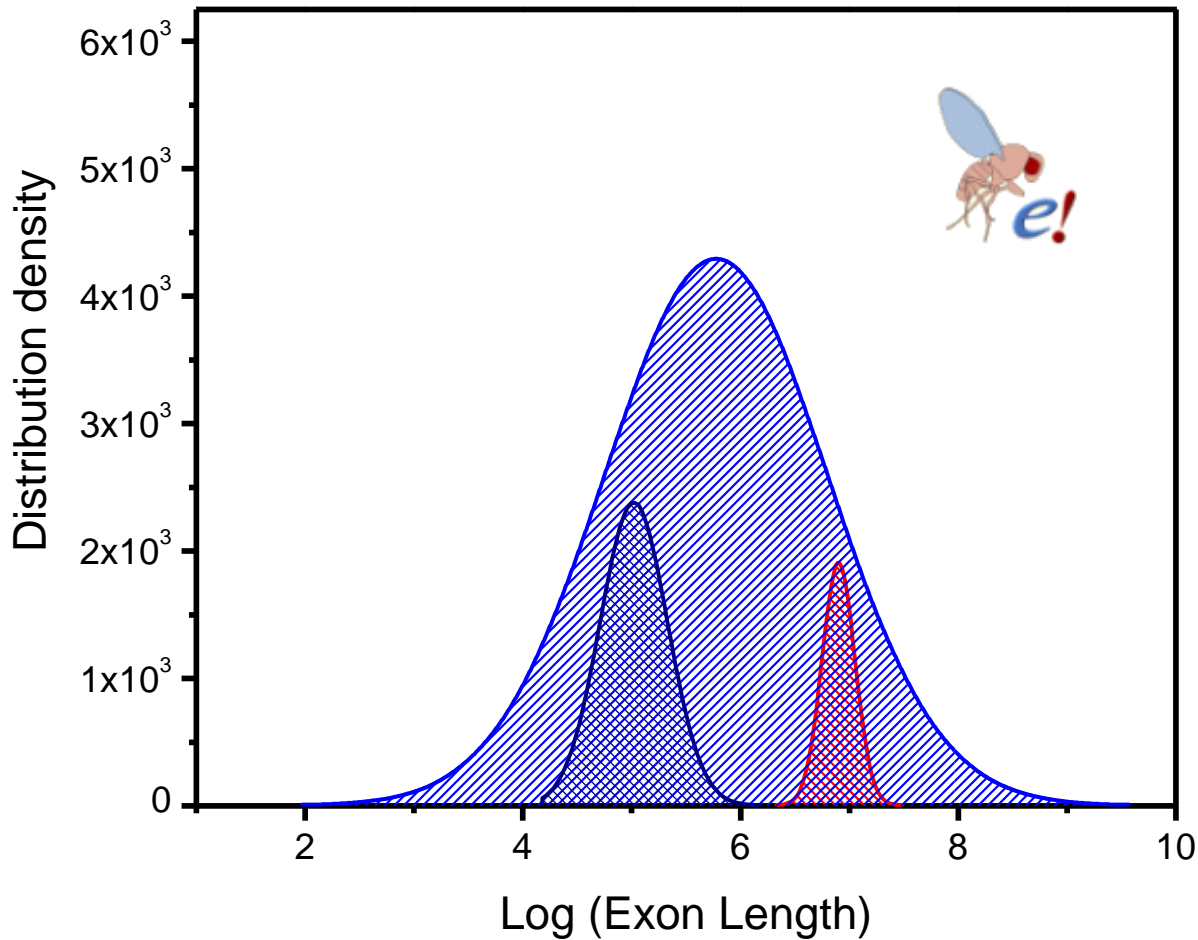
# Modeling exon size distributions in real genomes

50 000 time steps for  $p=0.001$



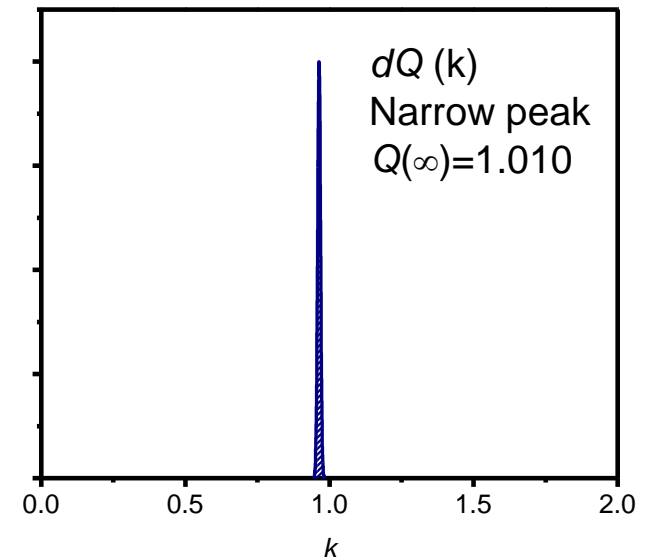
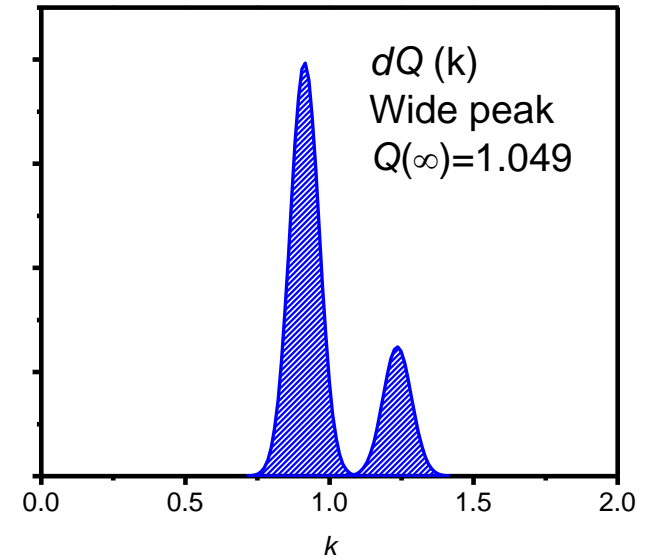
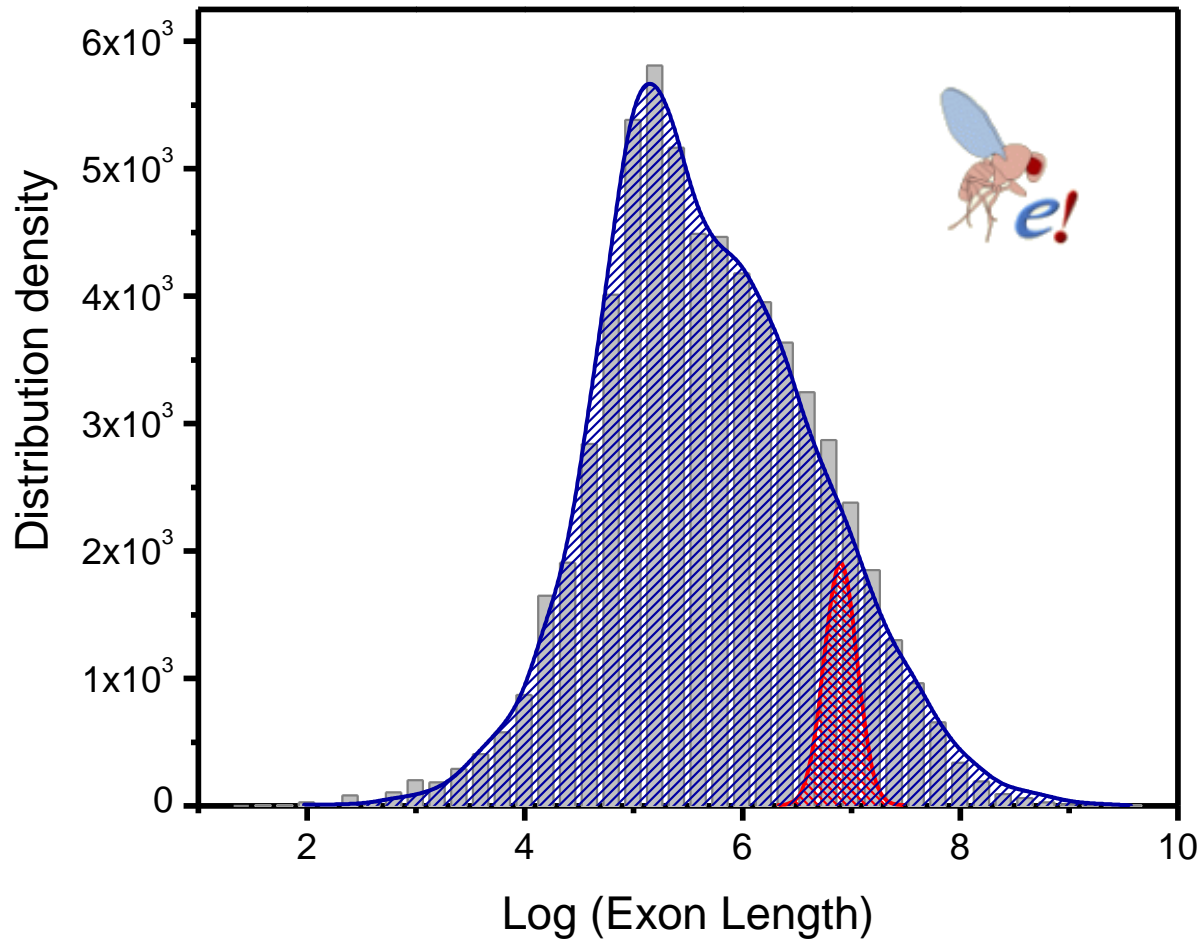
# Modeling exon size distributions in real genomes

50 000 time steps for  $p=0.001$

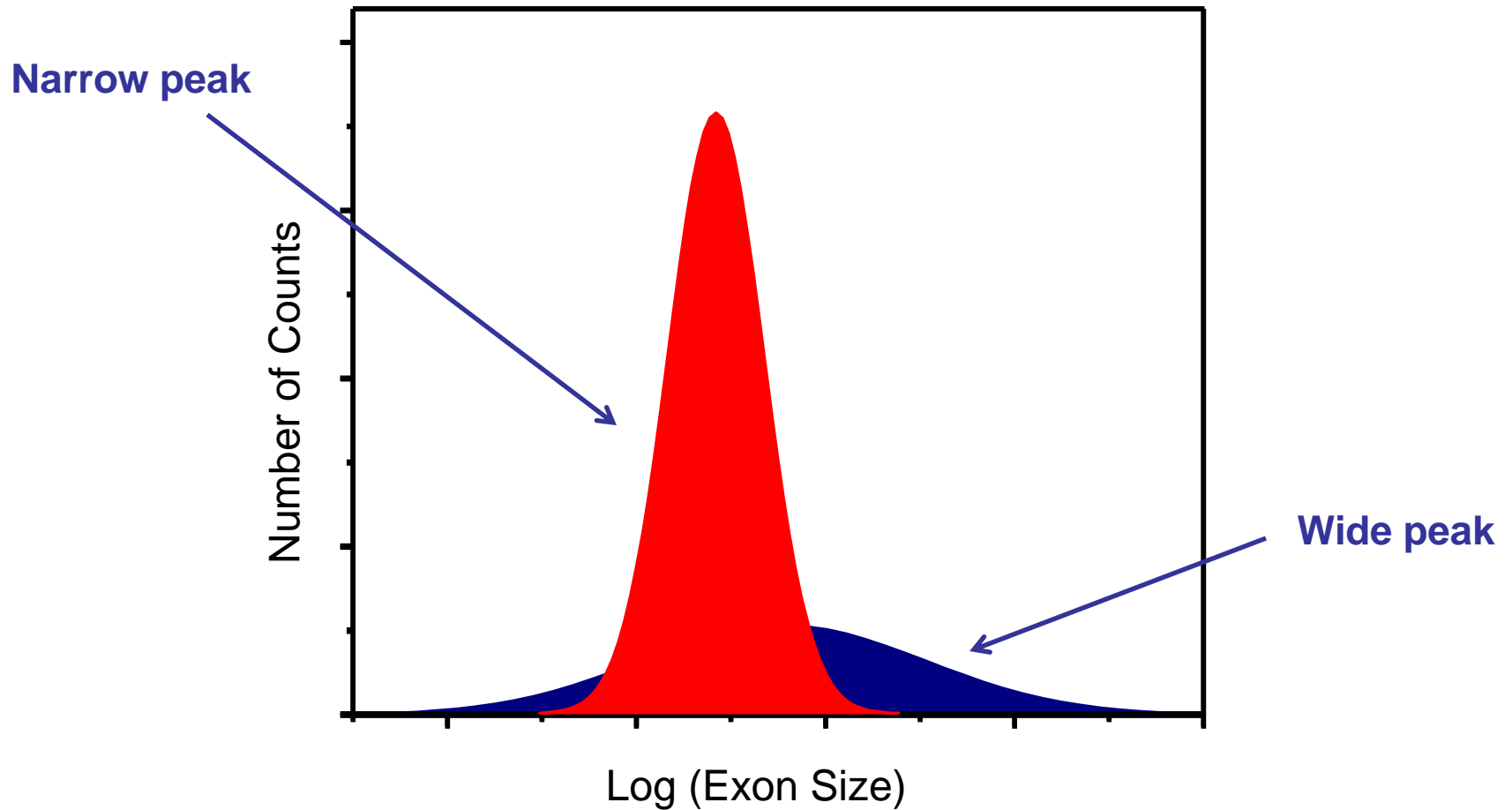


# Modeling exon size distributions in real genomes

50 000 time steps for  $p=0.001$

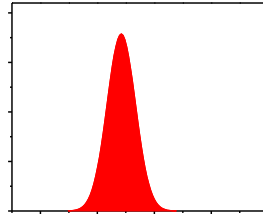


What could be the biological reason for two exons peaks?



## What could be the biological reason for two exons peaks?

### Narrow peak

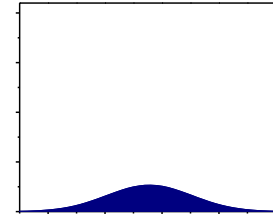


Holds approximately  
Constant position

Has approximately  
Constant width

Has greater relative  
occupation for  
Complex organisms

### Wide peak



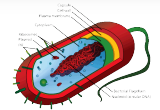
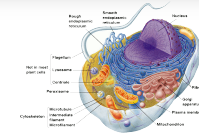
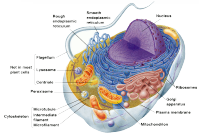
Changes position

Width increases with  
increasing complexity  
of and organism

Has greater relative  
occupation for  
Simple organisms



# Two ways of Gene Expression Regulation



**Alternative Splicing**

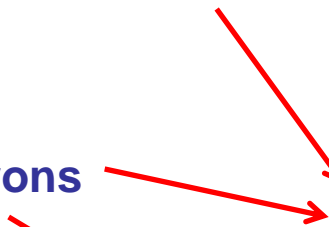
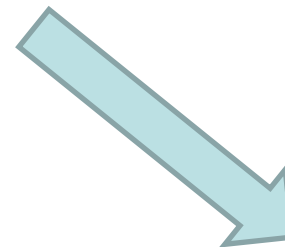
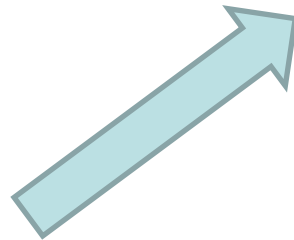
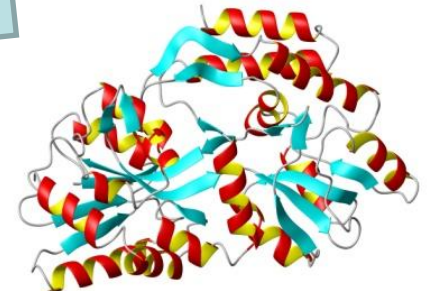
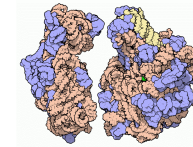
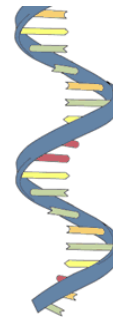
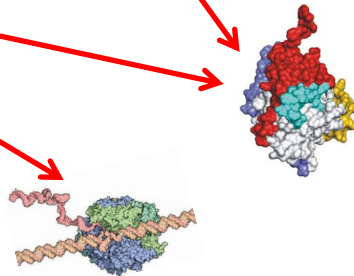
**mRNA**

**Untranslated regions of mRNA**

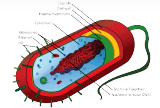
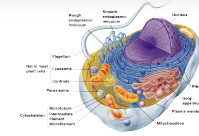
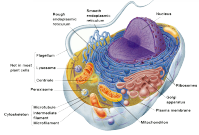
**Introns**

**DNA**

**Protein**



# Two ways of Gene Expression Regulation

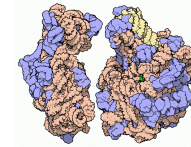
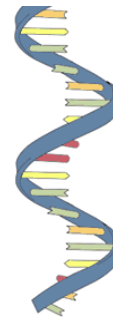


mRNA

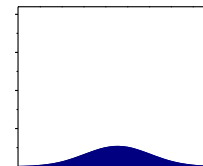
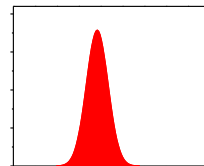
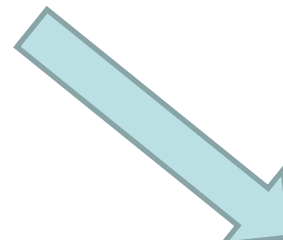
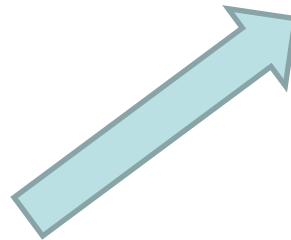
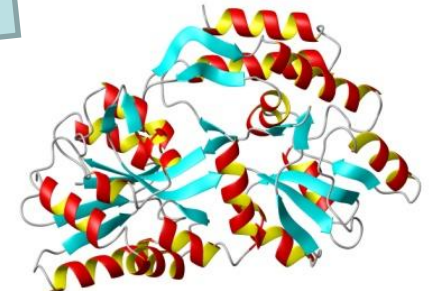
Untranslated  
regions of mRNA

Introns

DNA



Protein



# Conclusions

- Analysis of global properties of eukaryotic genomes reveals two distinct peaks in statistical distribution of exon sizes
- The observed peak could be fitted by a sum of two lognormal distributions which may imply that they originated by two different exon splitting pathways described in the general frameworks of Kolmogoroff splitting process
- Two observed peaks of exons could be correlated with the phenomenon of alternative splicing and with exons contributing into untranslated regions of mRNA. This suggests that the observed separation of exons in two different classes may be originated from two different ways of protein expression regulation.

# **Acknowledgments**

Michael Gribskov

Alexander Berezhkovskii