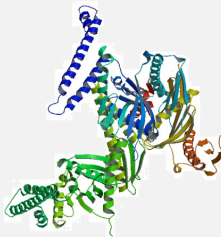


Analysis of translation termination sites in prokaryotes

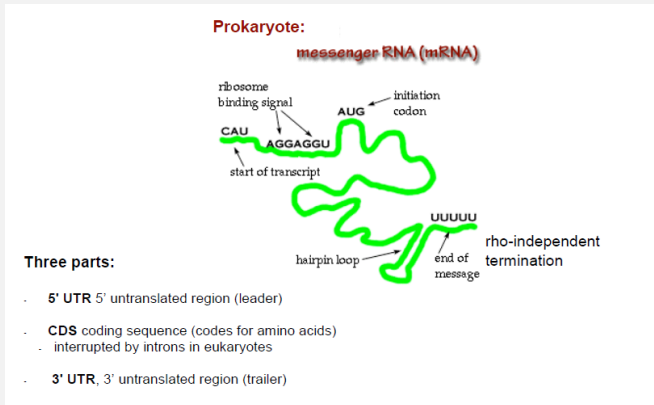
Przemysław Biecek, Paweł Mackiewicz, Dorota Mackiewicz,
Joanna Kiraga, Stanisław Cebrat

Wrocław/Warsaw University



- Basic facts about translation
- Statistical background
- Our data set
- Results for
 - DNA composition, positions from -50 to -1,
 - codon composition, positions from -30 to -1,
 - codon composition before different stop codons,
 - codon composition vs. different GC content.

General dogma: DNA \rightarrow mRNA \rightarrow proteins.



Initiation

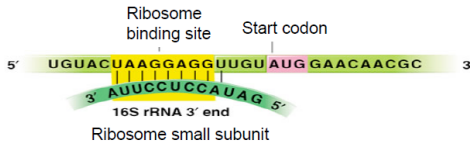
I. Initiation

Ribosome binds to Shine-Dalgarno sequence (AGGAGG).

Ribosome binding sites

complement to region in small subunit of ribosome

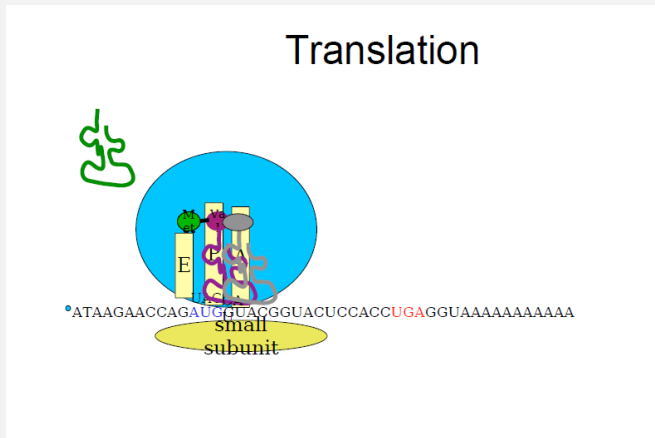
	Binding site sequences					Initiation codon				
Phage R17 A protein	UCC	UAG	GAG	GUU	UGA	CCU	AUG	CGA	GCU	UUU
Phage Q β replicase	UAA	CUA	AGG	AUG	AAA	UGC	AUG	UCU	AAG	ACA
Phage λ Cro	AUG	UAC	UAA	GGA	GGU	UGU	AUG	GAA	CAA	CGC
Phage Φ X174 A	AAU	CUU	GGA	GGC	UUU	UUU	AUG	GUU	CGU	UCU
<i>E. coli trpB</i>	AUA	UUA	AGG	AAA	GGA	ACA	AUG	ACA	ACA	UUA
<i>E. coli lacZ</i>	UUC	ACA	CAG	GAA	ACA	GCU	AUG	ACC	AUG	AUU
<i>E. coli RNA polymerase β</i>	AGC	GAG	CUG	AGG	AAC	CCU	AUG	GUU	UAC	UCC



From J. Russell, *Genetics*. Copyright © Pearson Education, Inc., publishing as Benjamin Cummings.

II. Elongation

tRNA with appropriate anticodon go to the ribosome site A.



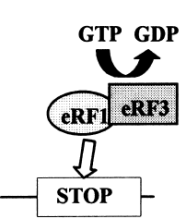
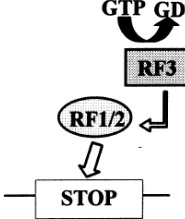
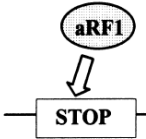
III. Termination

- Some release factor binds to stop codons.
- Peptide chain is released.
- In some cases ribosome start to translate next genes in other is dissociates.



Termination

- Stop codons may be different for different organisms (rare exceptions in small genomes).
- Release factors are different in different kingdoms.

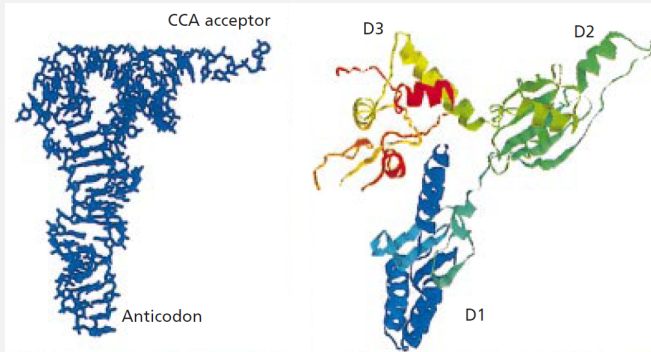
<i>Eukaryotes</i>	<i>Eubacteria</i>	<i>Archaea</i>
 <p><u>Codons recognised</u></p> <ul style="list-style-type: none"> • UAA, UAG, UGA • <i>Tetrahymena</i>; UGA • <i>Euplotes</i>; UAA, UAG 	 <p><u>Codons recognised</u></p> <ul style="list-style-type: none"> • RF1; UAA, UAG • RF2; UAA, UGA 	 <p><u>Codons recognised</u></p> <ul style="list-style-type: none"> • UAA, UAG, UGA

Genetic code

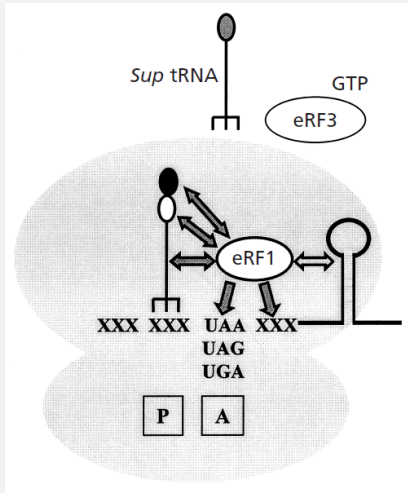
		2nd base			
		U	C	A	G
1st base	U	UUU (Phe/F)Phenylalanine	UCU (Ser/S)Serine	UAU (Tyr/Y)Tyrosine	UGU (Cys/C)Cysteine
		UUC (Phe/F)Phenylalanine	UCC (Ser/S)Serine	UAC (Tyr/Y)Tyrosine	UGC (Cys/C)Cysteine
		UUA (Leu/L)Leucine	UCA (Ser/S)Serine	UAA Ochre (Stop)	UGA Opal (Stop)
		UUG (Leu/L)Leucine	UCG (Ser/S)Serine	UAG Amber (Stop)	UGG (Trp/W)Tryptophan
	C	CUU (Leu/L)Leucine	CCU (Pro/P)Proline	CAU (His/H)Histidine	CGU (Arg/R)Arginine
		CUC (Leu/L)Leucine	CCC (Pro/P)Proline	CAC (His/H)Histidine	CGC (Arg/R)Arginine
		CUA (Leu/L)Leucine	CCA (Pro/P)Proline	CAA (Gln/Q)Glutamine	CGA (Arg/R)Arginine
		CUG (Leu/L)Leucine	CCG (Pro/P)Proline	CAG (Gln/Q)Glutamine	CGG (Arg/R)Arginine
	A	AUU (Ile/I)Isoleucine	ACU (Thr/T)Threonine	AAU (Asn/N)Asparagine	AGU (Ser/S)Serine
		AUC (Ile/I)Isoleucine	ACC (Thr/T)Threonine	AAC (Asn/N)Asparagine	AGC (Ser/S)Serine
		AUA (Ile/I)Isoleucine	ACA (Thr/T)Threonine	AAA (Lys/K)Lysine	AGA (Arg/R)Arginine
		AUG (Met/M)Methionine, Start ¹¹	ACG (Thr/T)Threonine	AAG (Lys/K)Lysine	AGG (Arg/R)Arginine
	G	GUU (Val/V)Valine	GCU (Ala/A)Alanine	GAU (Asp/D)Aspartic acid	GGU (Gly/G)Glycine
		GUC (Val/V)Valine	GCC (Ala/A)Alanine	GAC (Asp/D)Aspartic acid	GGC (Gly/G)Glycine
		GUA (Val/V)Valine	GCA (Ala/A)Alanine	GAA (Glu/E)Glutamic acid	GGA (Gly/G)Glycine
		GUG (Val/V)Valine	GCG (Ala/A)Alanine	GAG (Glu/E)Glutamic acid	GGG (Gly/G)Glycine

Release Factors

Release factors (1 and 2) are structurally similar to tRNA.



Context of stop codons



Context of the termination site

- last few amino acids,
- nucleotides before stop,
- nucleotides after stop,
- mRNA structure (hairpins, regions rich in nucleotide A).

- There is a lot of papers about stop translation signal. Usually only one organism is analyzed e.g. *E. coli* or other model organisms.
- RF speciation
 - RF1 recognise codons TAG and TAA,
 - RF2 recognise codons TGA and TAA.
- RF effectiveness
 - TAA (ochre) — the most frequent, strong signal, fails one per 1000 passes,
 - TGA (opal) — a weak signal of termination and fails one per 4 passes.
- Weak stops are natural source of alternative protein products!
- Tandem Stop Codons, higher than expected number of codon TAA in position +3 after stop.
TSC are more frequent in highly expressed genes.

Comparative study for over 400 prokaryotic genomes (all prokaryotic genomes completed before VI 2007, $\approx 1\,000\,000$ genes).

- Goal 1: To identify signals correlated with stop translation.
- Goal 2: To identify factors which may explain evolution of such signals.

Data set from NCBI

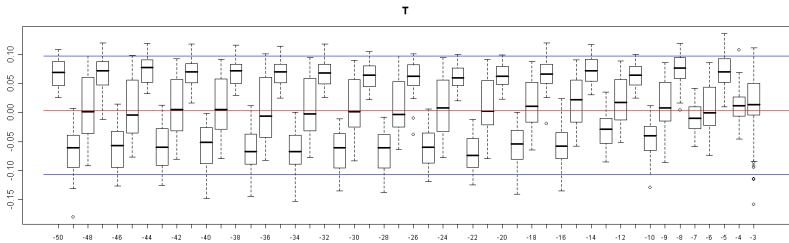
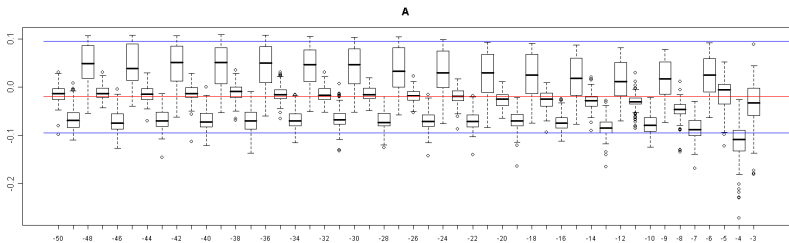
- genes from plastids and from chromosome are analyzed together,
- one candidate is chosen for a set of genomes with identical fourth-level taxonomic classification,
- .fna files with gene sequences and .ptt files with gene coordinates.

For every genome the composition skew is computed

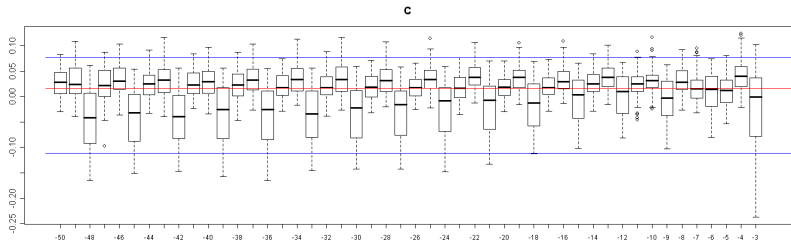
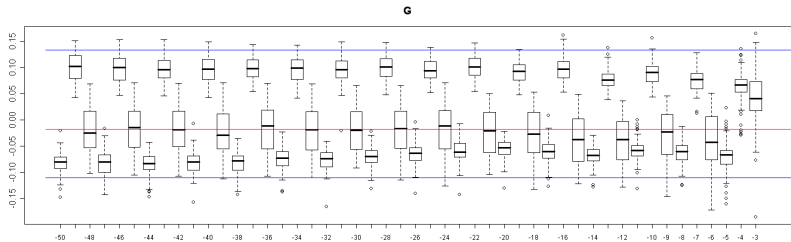
$$skew = \hat{\rho}_{-i,x} - \hat{\rho}_{intra,x},$$

where $\hat{\rho}_{-i,x}$ is the codon/nucleotide x frequency at position $-i$ before stop, while $\hat{\rho}_{intra,x}$ is the intragene codon/nucleotide x frequency.

DNA rythm



DNA rythm



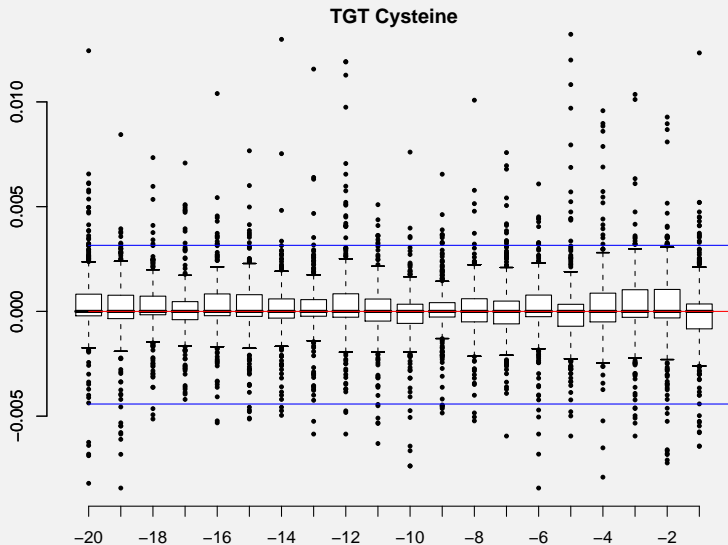
Changes in nucleotide composition

- loss of tri-nucleotide frequency pattern,
- nucleotide A is overrepresented,

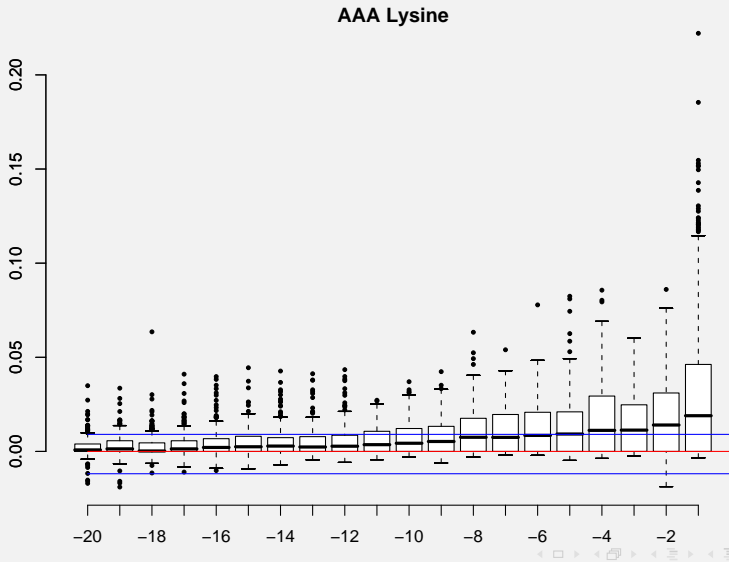
have impact on mRNA secondary structure and may prepare ribosome to termination.

We think that codons composition is more important than nucleotide composition.

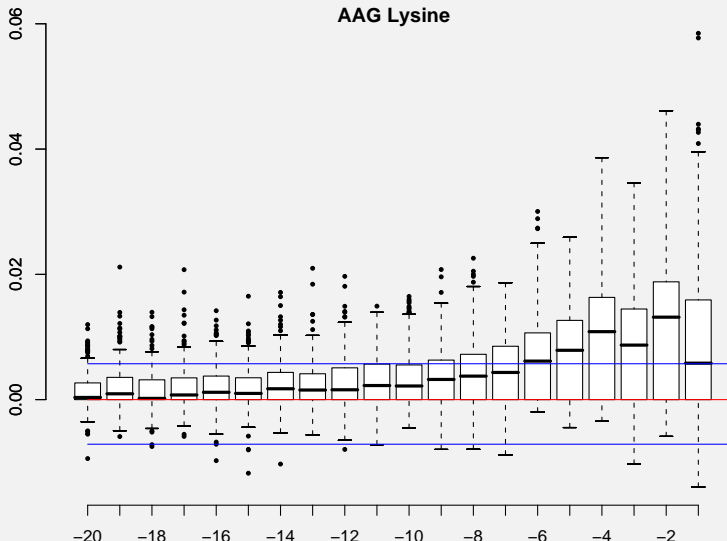
Codon composition



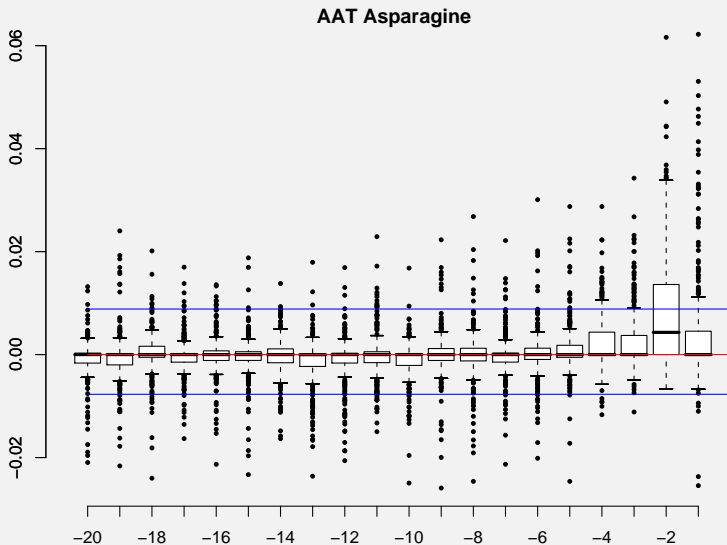
Codon composition



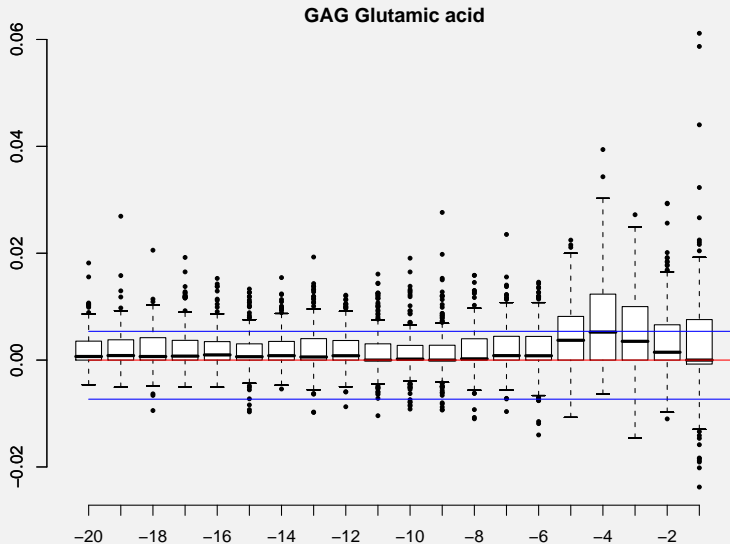
Codon composition



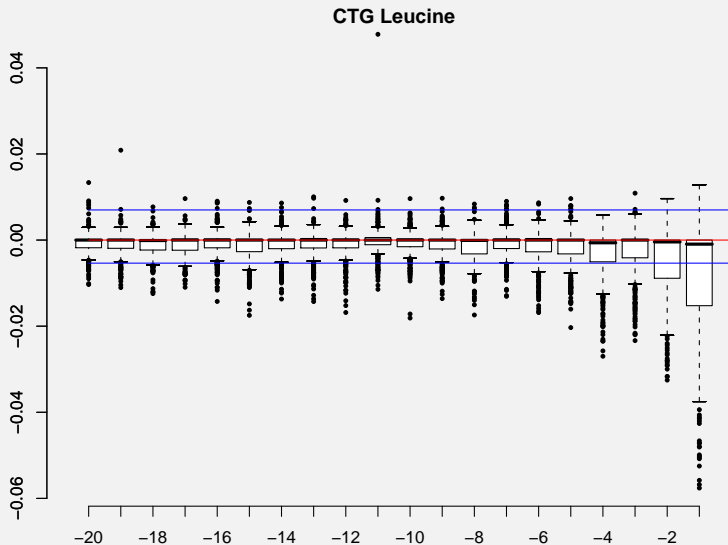
Codon composition



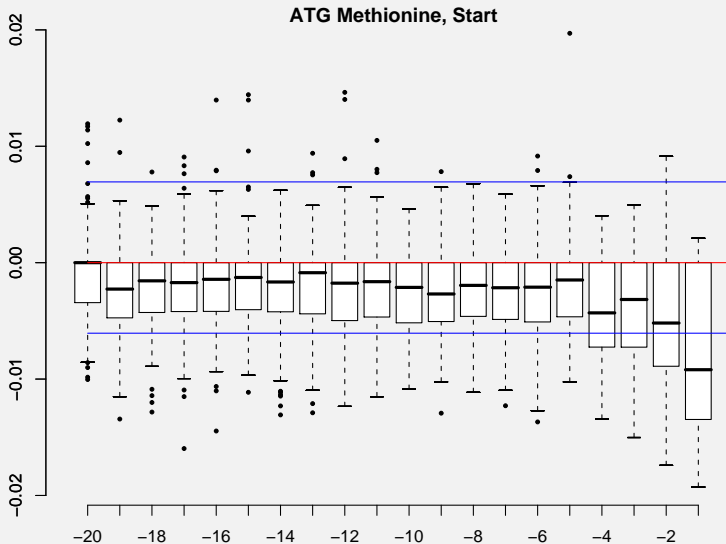
Codon composition



Codon composition



Codon composition



Changes in codons composition

- codon AAA is highly overrepresented,
- in some codons (e.g. AAG) strong signal is observed in more than one position,
- in most cases the strongest signal is in position -1,
- in some codons (e.g. AAT) in position -2 or (e.g. GAG) in position -4,
- some codons (e.g. ATG - start) are strongly underrepresented.

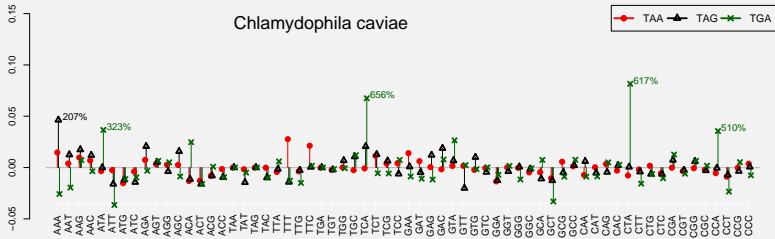
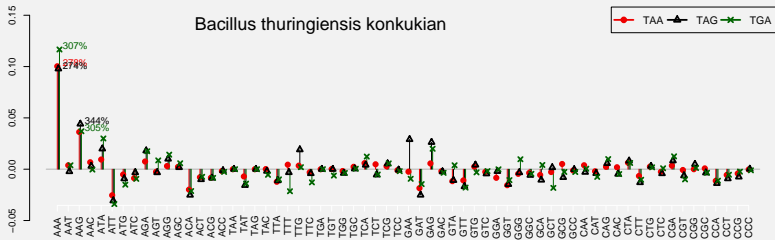
Let's have a look on some genomes. We plot results for position -1.
As before the codon skew is computed as:

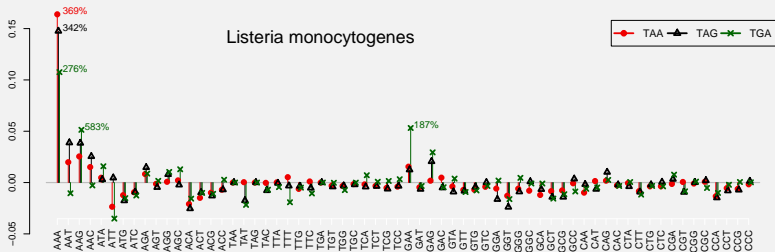
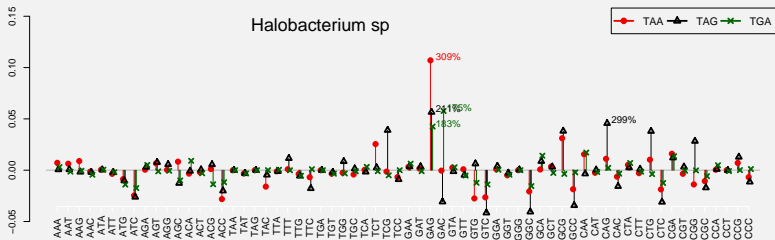
$$skew_x = \hat{\rho}_{-i,x} - \hat{\rho}_{intra,x},$$

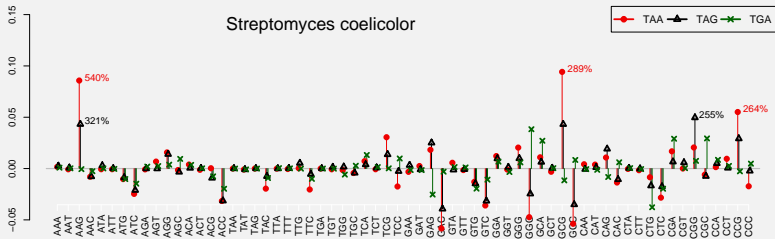
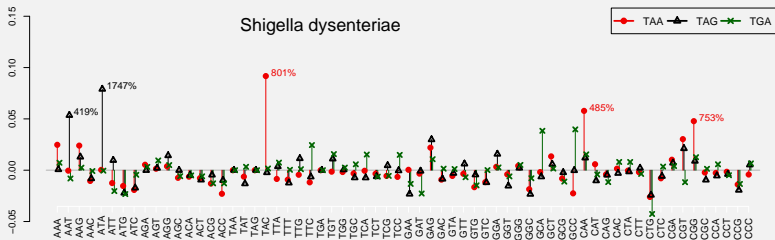
codon relative change is computed as

$$relative_x = \hat{\rho}_{-i,x} / \hat{\rho}_{intra,x},$$

where $\hat{\rho}_{-i,x}$ is the codon x frequency at position $-i$ before stop,
while $\hat{\rho}_{intra,x}$ is the codon x intragene frequency.







In some genomes (e.g. *Bacillus thuringiensis*) there is no significant changes in codon overrepresentation before different stop codons.

In some genomes (e.g. *Shigella dysenteriae*) there is significant difference.

Why they behave differently?

The test of proportions was used to compare the frequency of given codon in position before a stop signal and in positions inside genes (intragenic). The frequency before stop signal was also compared with intergenic frequency, and in most cases results are similar.

The Bonferroni correction was applied to deal with the number of tests (for 400 genomes and 64 codons). For a particular test the very conservative significance level was used

$$\alpha = \frac{0.05}{400 * 64}.$$

While we test on such small significance level, all rejections are positive with high probability.

Genomes were divided into 4 groups according to their GC content.

GC content	<37.5	37.5÷50	50÷62.5	>62.5
AAA	90.57%	82.68%	80.61%	73.42%
AAC	23.58%	18.90%	9.18%	7.59%
AAG	57.54%	49.52%	42.85%	29.11%
AGA	26.42%	51.18%	66.33%	75.95%
AGG	27.36%	24.41%	32.65%	54.43%
TCA	5.66%	19.69%	65.31%	92.41%
GGA	0.94%	5.51%	44.90%	86.08%
GCA	0.94%	6.30%	56.12%	96.20%
CGA	14.15%	26.77%	90.82%	100.00%
CCA	0.94%	3.94%	47.96%	94.94%

- Some codons (e.g. GCA, CGA, TCA) are often overrepresented in position -1 in genomes with high GC content. Other codons (e.g. AAA, AAG) are often overrepresented in position -1 in genomes with low GC content! (no artefacts)
- The lower overrepresentation of AAA codon in genomes with high GC content may result from different mutational pressure.

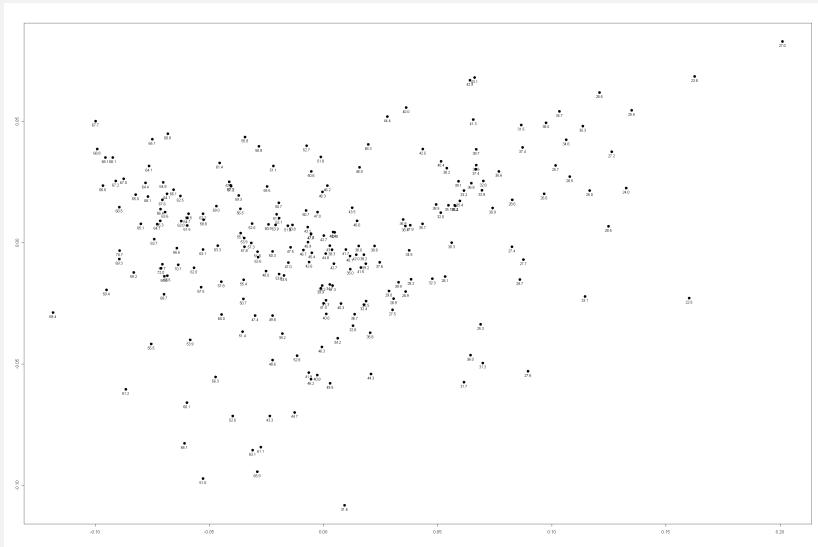
Note that codon frequencies before stop were compared with intergenic or intragenic codon frequencies in the same genome.

We use Multidimensional Scaling to visualize differences in codon composition for different genomes.

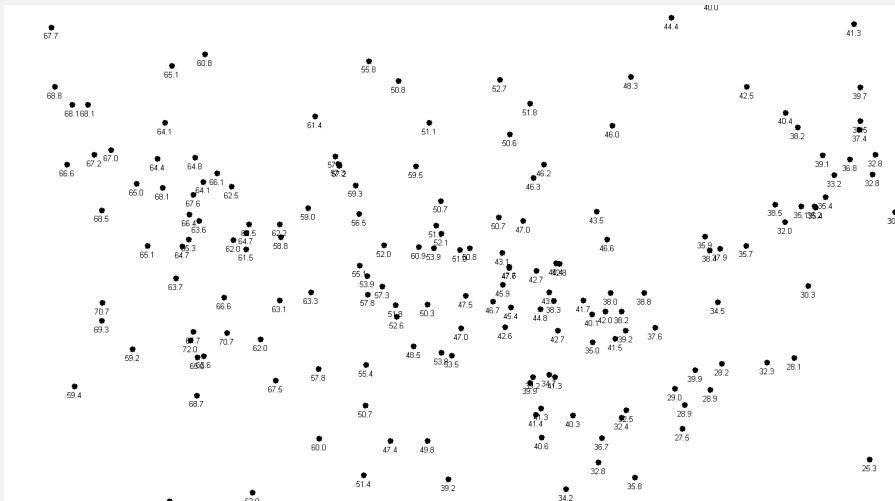
MDS is a set of related statistical techniques often used in data visualization or feature extraction problems.

Using MDS we find a new coordinates for genomes in 2D or 3D space. Coordinates are computed in a way to minimize differences between distances in the original parameter space (in our case it has 128 dimensions) and in the new parameter space.

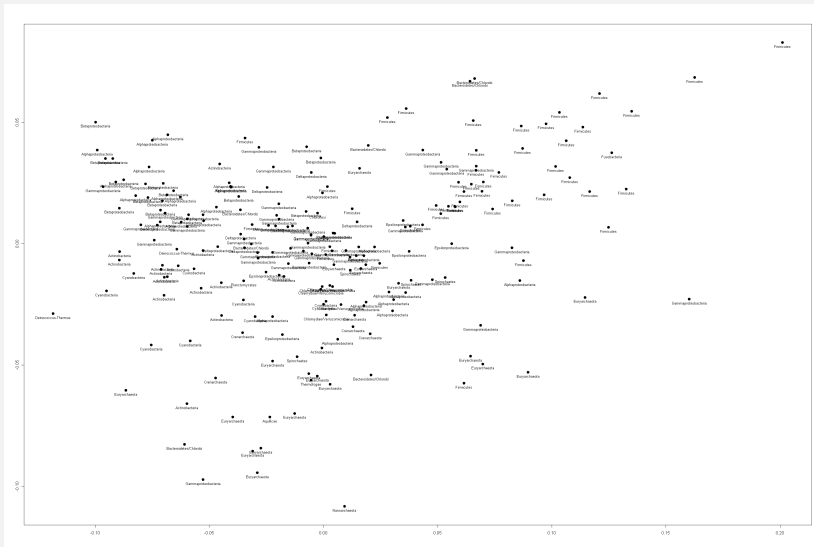
MDS for GC content



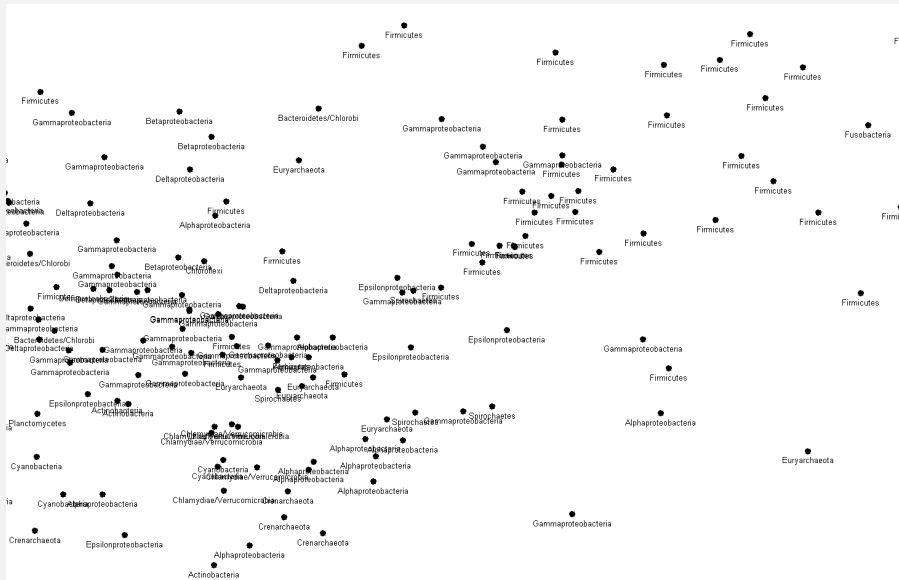
MDS for GC content



MDS for taxonomic groups



MDS for taxonomic groups

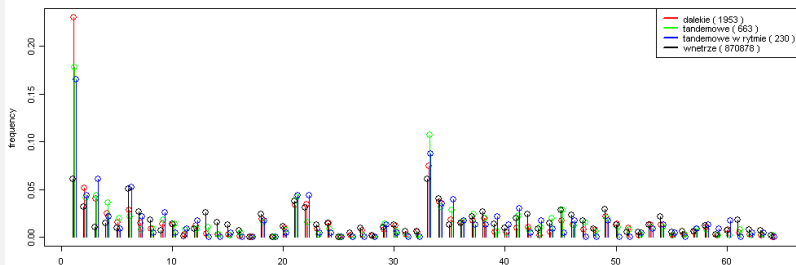
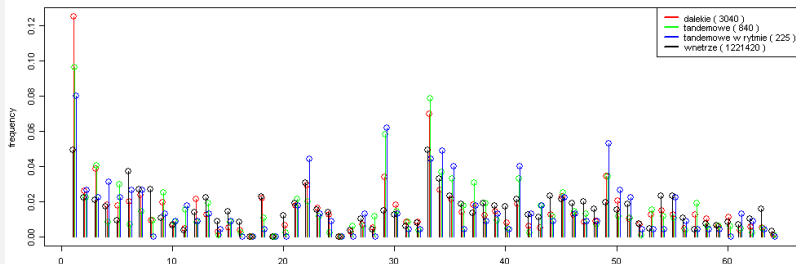


There is much more figures that may be presented

- results for Tandem genes,
- other genome properties (size, habitat, temp. range, etc.),
- results for amino acids composition,
- results for different isoelectric points (pI),
- ...

but they are not so interesting.

Other results



- Nucleotide and amino acid composition is changed before stop codon.
- Codons **AAA**, **AGA** and **AAG** are highly overrepresented.
- Codons **GAA**, **GGA** and **GAG** are overrepresented.
- Codon AUG is underrepresented.
- Differences in codon overrepresentation correlates with GC content and (possible) mutation pressure.