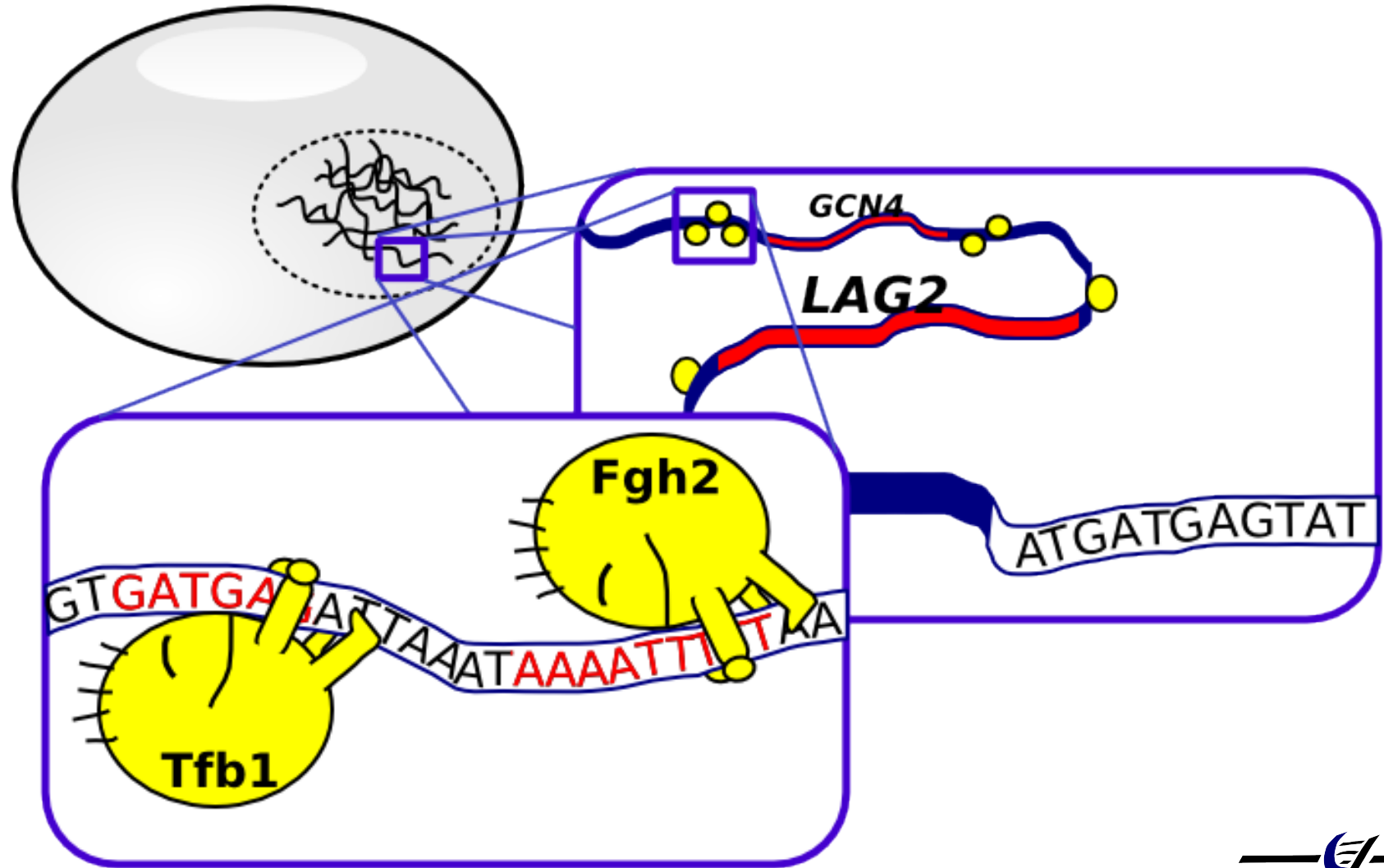

An Evolutionary Model of DNA Substring Distribution

Konstantin Tretyakov (kt@ut.ee)
(joint work with Meelis Kull and Jaak Vilo)

University of Tartu
Bioinformatics & Data Mining Group

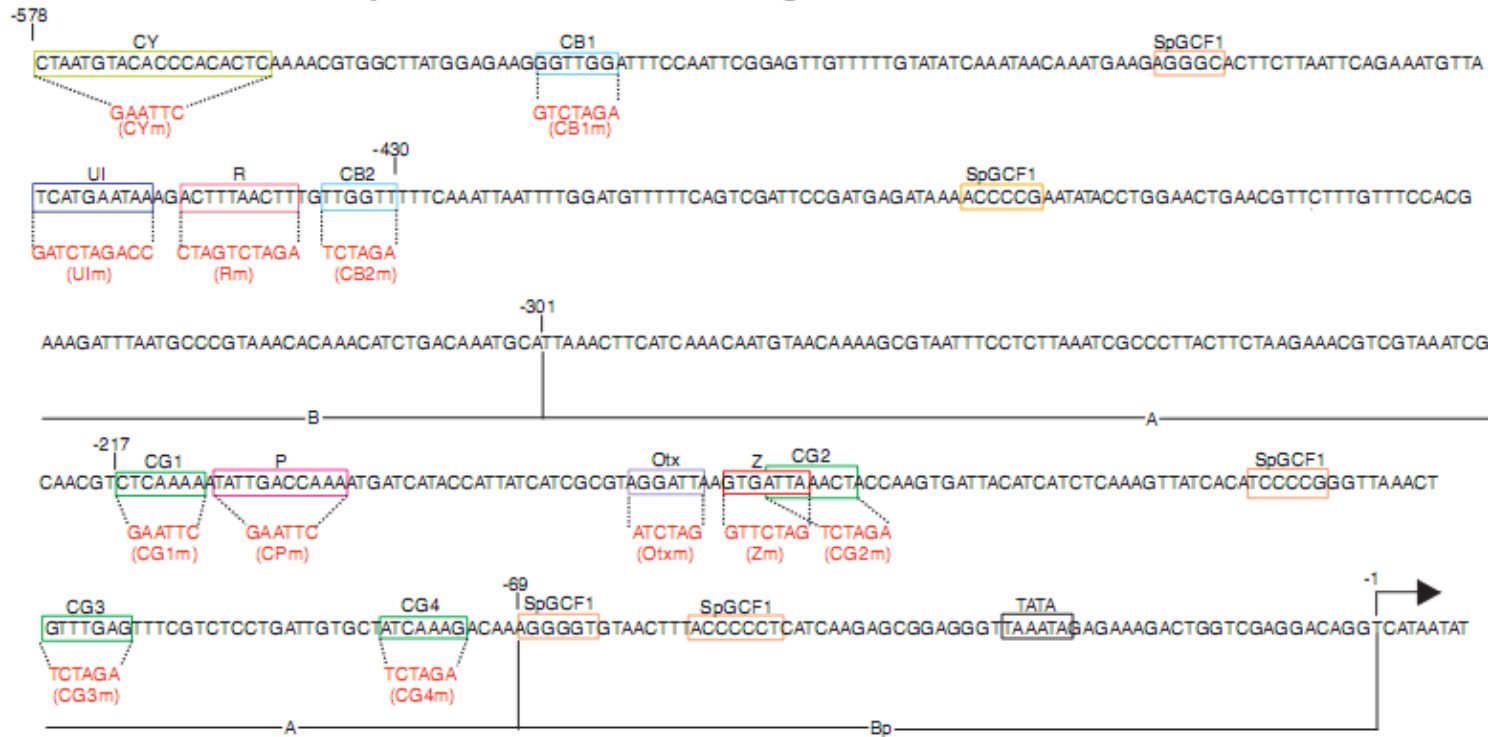


DNA Regulatory Regions



DNA Regulatory Regions

Promoter sequence of Endo 16 gene



Purple sea urchin



Motif Discovery

FGH: ATCG**GTGATGAGTT**AGACG
G**GAGATGCGTT**ATCATGTAG
AGATCGAGATACGATAGBAA
GAAAATAGC**GAGATGAGTA**

HGH: TATAATCG**GTGATGAGTT**AG
AGAGAGABACGGAGATACGA
TAT**CAGATCAGATT**GGTA

Modeling "Background"

Sequence:

CTAATGTACACCCACACTCAAACGTGGCTTATGGAGAAGGGTTGGATTTCCAATT
ACTTCTTAATTCAGAAATGTTATCATGAATAAAGACTTTAACTTTGTTGGTTTTTC
AAAACCCCGAATATACCTGGAAGTGAACGTTCTTTGTTTCCACGAAAGATTTAATG
CAAACAATGTAACAAAAGCGTAATTTCTCTTAAATCGCCCTTACTTCTAAGAAAC
TCATACCATTATCATCGCGTAGGATTAAGTGATTAAACTACCAAGTGATTACATCA
TCGTCTCCTGATTGTGCTATCAAAGACAAAGGGGTGTAACTTTACCCCTCATCAA

Background:

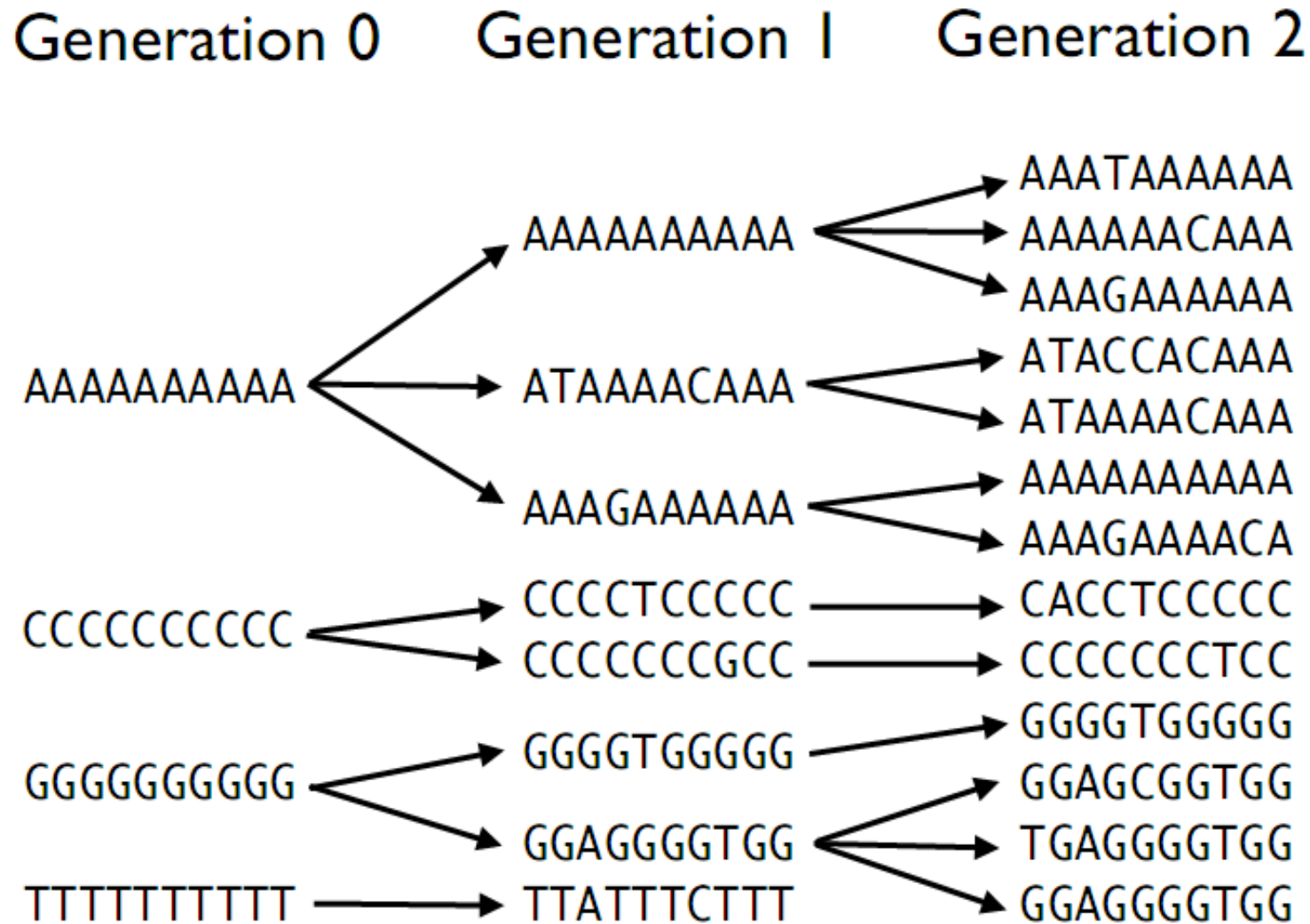
AA	- 14%	GA	- 6%
AC	- 6%	GC	- 2%
A	- 34%	AG	- 5%
C	- 19%	AT	- 9%
G	- 18%	CA	- 7%
T	- 30%	TA	- 7%
		CC	- 4%
		TC	- 7%
		CG	- 3%
		TG	- 6%



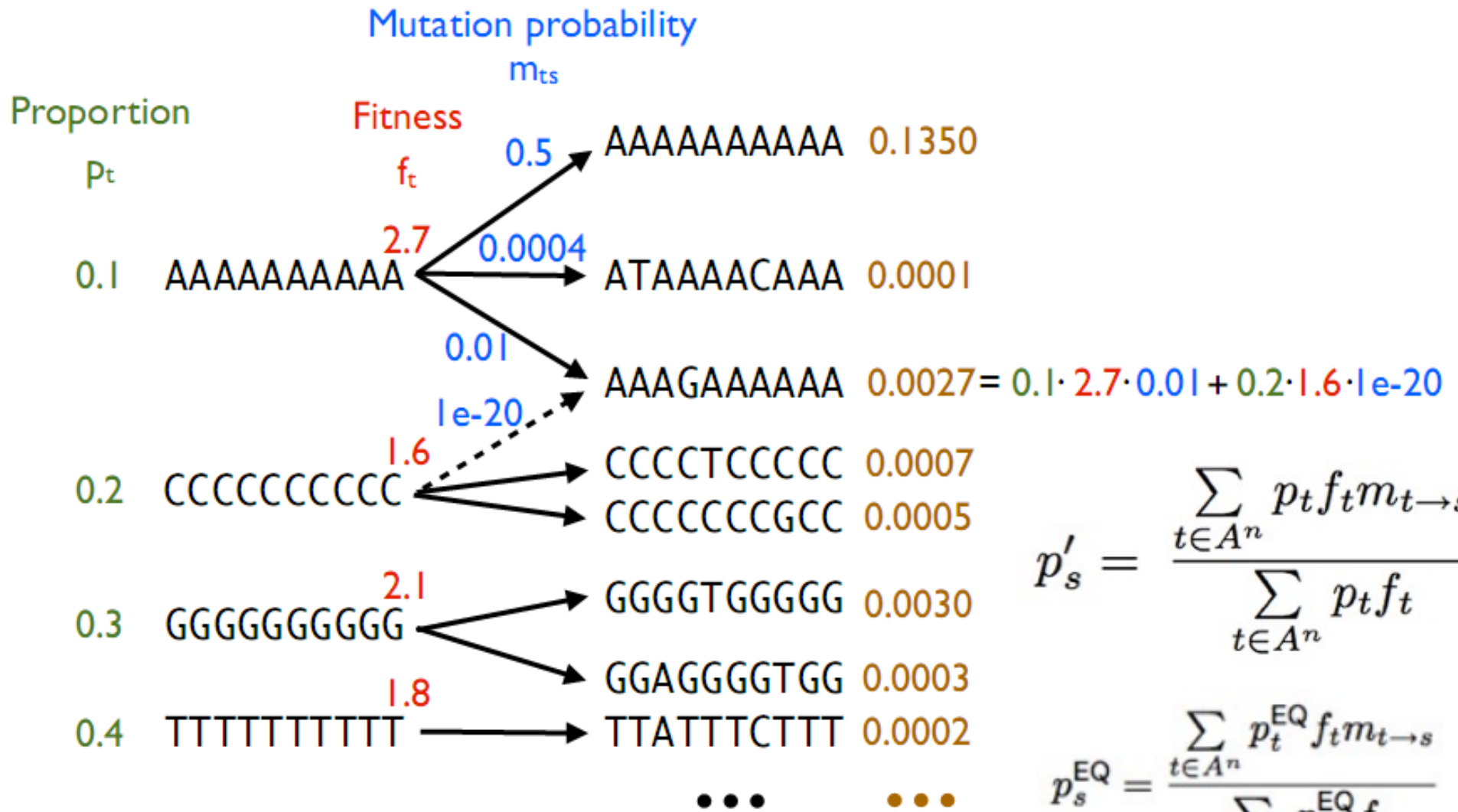
Modeling "Background"

- How does a DNA substring distribution **arise** in the first place?
 - Evolutionary model

Evolutionary Model



Evolutionary Model



$$p'_s = \frac{\sum_{t \in A^n} p_t f_t m_{t \rightarrow s}}{\sum_{t \in A^n} p_t f_t}$$

$$p_s^{EQ} = \frac{\sum_{t \in A^n} p_t^{EQ} f_t m_{t \rightarrow s}}{\sum_{t \in A^n} p_t^{EQ} f_t}$$

k-mer equilibrium distribution

	p_s^{EQ}	$k=3$	
AAAAAAAAAA	1.5e-5	AAA	0.052
AAAAAAAAAC	4.3e-6	AAC	0.027
AAAAAAAAAG	3.6e-6	AAG	0.023
AAAAAAAAAT	5.2e-6	AAT	0.015
AAAAAAAA <u>ACA</u>	7.1e-6	ACA	0.043 = 1/10 · 7.1e-6 + 2/10 · 3.5e-5 + ...
...		ACC	0.018
		ACG	0.012
		ACT	0.010
CAGACAGCAA	3.5e-5	AGT	0.016
<u> </u> <u> </u> <u> </u> <u> </u>		ACA	0.013
...		...	
TTTTTTTTTC	9.8e-6	TTC	0.021
TTTTTTTTTG	1.0e-5	TTG	0.027
TTTTTTTTTT	2.6e-5	TTT	0.033

Main contribution

- **Given the fitness function and mutation rates, compute the expected frequencies of K-mers.**

- Mustonen and Lässig. Evolutionary population genetics of promoters: predicting binding sites and functional phylogenies. *Proc Natl Acad Sci USA* (2005) vol. 102 (44) pp. 15936-41
- Moses et al. Large-scale turnover of functional transcription factor binding sites in *Drosophila*. *PLoS Comput. Biol.* (2006) vol. 2 (10) pp. e130
- Huang et al. Phylogenetic simulation of promoter evolution: estimation and modeling of binding site turnover events and assessment of their impact on alignment tools. *Genome Biol* (2007) vol. 8 (10) pp. R225
- Doniger and Fay. Frequent gain and loss of functional transcription factor binding sites. *PLoS Comput. Biol.* (2007) vol. 3 (5) pp. e99



k · mer distribution

$$p_a^{\text{EQ}} = \frac{\sum_{b \in A^k} m_{b \rightarrow a} \sum_{u \in A^{n-k}} p_{b \cdot u}^{\text{EQ}} f_{b \cdot u}}{\sum_{b \in A^k} \sum_{u \in A^{n-k}} p_{b \cdot u}^{\text{EQ}} f_{b \cdot u}}$$

$$\sum_{u \in A^{n-k}} p_{b \cdot u}^{\text{EQ}} f_{b \cdot u} \approx \frac{p_b^{\text{EQ}}}{|A^{n-k}|} \sum_{u \in A^{n-k}} f_{b \cdot u}$$

$$p_a^{\text{EQ}_k} = \frac{\sum_{b \in A^k} p_b^{\text{EQ}_k} \cdot f_b m_{b \rightarrow a}}{\sum_{b \in A^k} p_b^{\text{EQ}_k} f_b}$$



Experiments

- To check approximations we have compared p^{EQ} and p^{EQ_k}
- Calculation of p^{EQ} is hard, thus we took sequences of length 8 of letters A and C
- Independent point mutations with probability r at each location
- Fitness defined by the number of occurrences of a motif q

Results

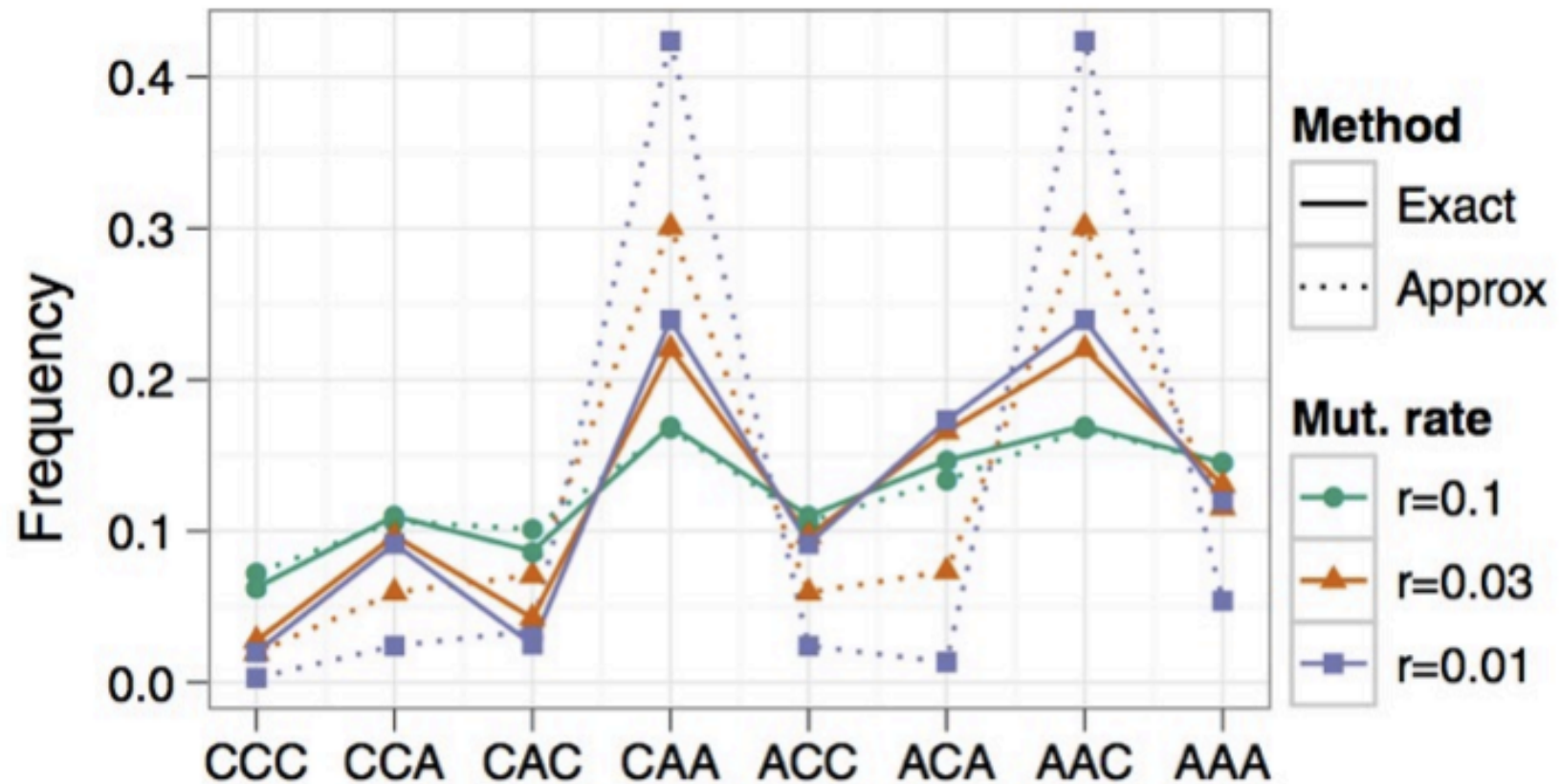


Fig. 2. The exact and approximate 3-mer distributions for point-mutation rates $r = 0.1, 0.03, 0.01$ where fitness is defined with strategy (S2) for substring $q = AAC$

Results

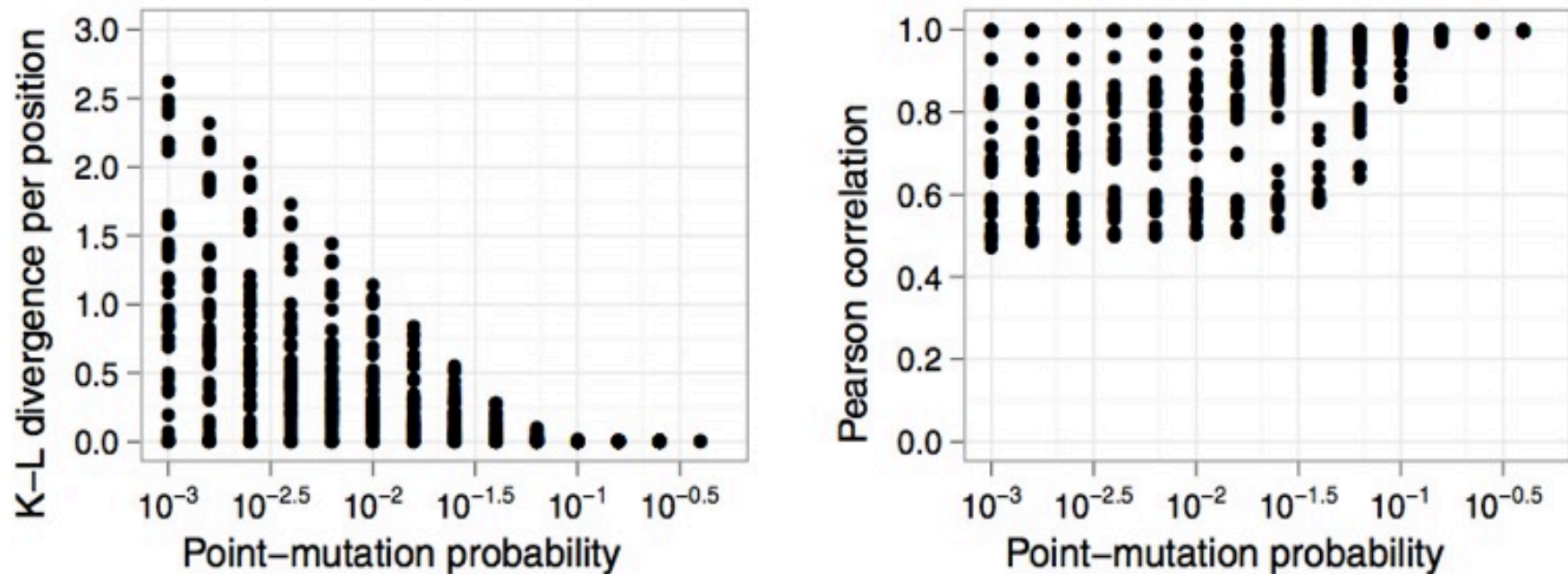


Fig. 1. The effect of point-mutation rate r on the approximation quality measured as Kullback-Leibler divergence per position and correlation between the exact and approximated distributions. Each circle denotes an experiment with a different set of parameters.

F u t u r e w o r k

- **Matching the model on data**
 - **Estimating model parameters from data**
 - **When and to what extent does the genetic equilibrium hold**
- **Improving approximation**
- **Applications to motif discovery**

Thanks

