

Monotony of Surprise and Large-Scale Quest for Unusual Words

Stefano Lonardi

University of California, Riverside

joint work with A. Apostolico, M. E. Bock, F. Gong

Detection of unusual words

- **GIVEN**
 - a text x
 - a probabilistic *model* of the source which has generated x
- **FIND** all the substrings of x which are significantly more *frequent/rare* than the model-based expectation

Example

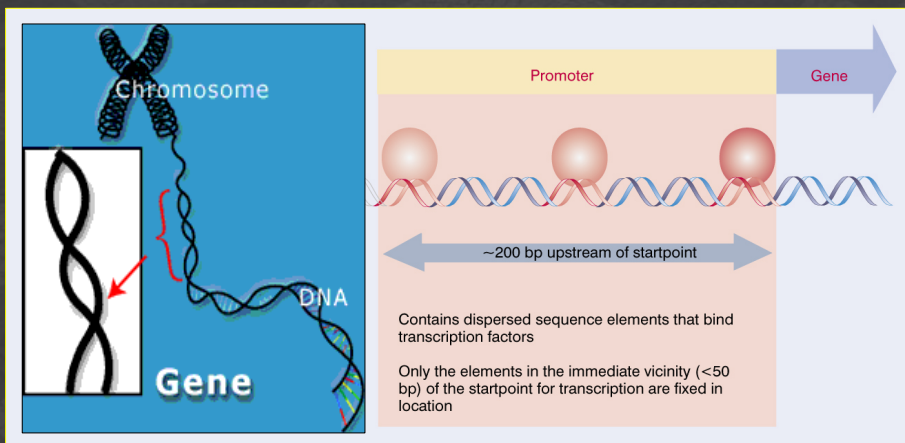
```
...ATGACAAGTCCTAAAAAGAGCGAAAACACAGGGTTGTTTGATTGTAGAAAATCACAGCG  
>MEK1  
CCACCCTTTTGTGGGGCTTCTATTTCAAGGACCTTCATTATGGAACAGGGCGAGGTTGT  
TTGTTCTTCTGTCATGTTGCGGCGAGTGCCTAAGAAAGCGGGACGTAAGCAGTTTAGCCA  
TTCTAAAAGGGGCATTATCAGAATAAGAAGGCCCTATGAGGTATGATTGTAAAGCAAGTG  
GTGTAAAATTGTGTGCTACCTACCGTATTAGTAGGAACAATTATGCAAGAGGGGTCCTGT  
GCAAAATAAAAAATATATATCTAGAAAAAGAGTAGGTAGGTCCTTCACAATATTGACTGAT  
AGCGATCTCCCTCACTATTTTCACTTATATGCAGTATATTTGTCTGCTTATCTTTCATTA  
AGTGAATCATTGTAGTTTATTCTACTTTATGGGTATTTTCCAATCATAAAGCATACC  
GTGTAATTTAGCCGGGAAAAGAAGAAATGATGGGGCTAAATTCGGGGC...
```

parameters

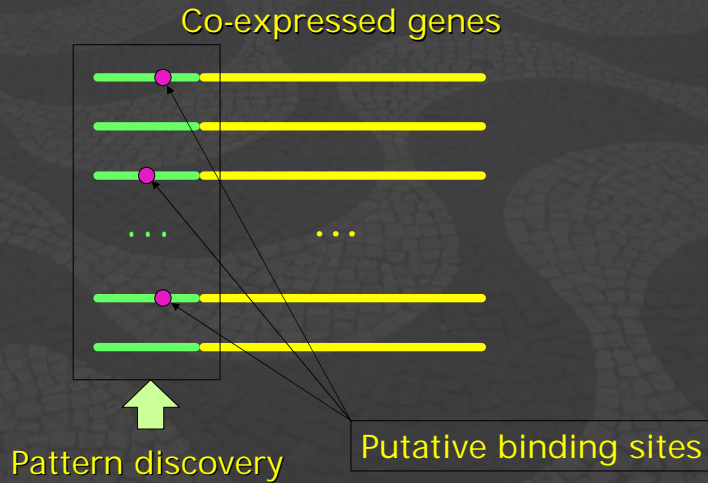
MODEL

?

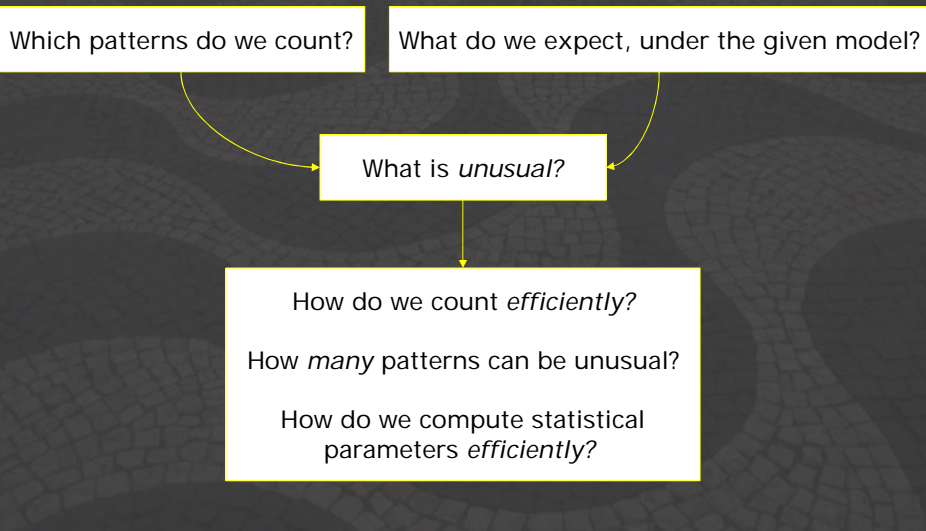
Transcription factors binding sites



Transcription factors binding sites



General framework



Notations

x : sequence, $|x| = n$

y : substring of x , $|y| = m$

$f(y)$: number of occurrences of y in x

Bernoulli model

Let Z_y be a r.v. for the number of occurrences of y ,
 p_a be the probability of $a \in \Sigma$, and $|y| = m \leq (n+1)/2$

- $E(Z_y) = (n - m + 1) \prod_{i=1}^m p_{y_{[i]}} = (n - m + 1) \hat{p}$
- $Var(Z_y) = E(Z_y)(1 - \hat{p}) - \hat{p}^2 (n - m + 1)(n - m) + 2\hat{p}B(y)$

where $B(y) = \sum_{d \in P(y)} (n - m + 1 - d) \prod_{i=m-d+1}^m p_{y_{[i]}}$

and $P(y)$ is the set of period lengths of y

Scores

$$z_1(y) = f(y) - E(Z_y)$$

$$z_2(y) = \frac{f(y) - E(Z_y)}{\sqrt{E(Z_y)}}$$

$$z_3(y) = \frac{f(y) - E(Z_y)}{\sqrt{E(Z_y) (1 - \hat{p})}}$$

$$z_4(y) = \frac{f(y) - E(Z_y)}{\sqrt{\text{Var}(Z_y)}}$$

where Z_y is a r.v. for the number of occurrences of y

What is "unusual" ?

Definition

Let y be a substring of x and $T \in \mathbb{R}^+$

- if $z(y) > T$, then y is **over-represented**
- if $z(y) < -T$, then y is **under-represented**
- if $|z(y)| > T$, then y is **unusual**

Problem definition

Given

- Sequence x
- Model M
- Type of count (f, \dots)
- Score function z
- Threshold T

Find

- The set of all unusual words in x
w.r.t. $(f/\dots, z, M, T)$

Computational problems

- Counting "events" in strings
(occurrences, ...)
- Computing expectations, variances,
and scores (under the given model)
- Detecting and visualizing unusual
words

Combinatorial problem

- A sequence of size n could have $O(n^2)$ unusual words
- How to limit the set of unusual words?

Monotony of surprise

Theorem

Let C be a subset of words from text x . If $f(y)$ remains **constant** for all y in C , then any score of the type

$$z(y) = \frac{f(y) - E(y)}{N(y)}$$

is monotonically **increasing** with $|y|$ provided that

- $N(y)$ is monotonically **decreasing** with $|y|$
- $E(y)/N(y)$ is monotonically **decreasing** with $|y|$

Theorem

Score functions

$$z(y) = f(y) - E(Z_y)$$

$$z(y) = \frac{f(y) - E(Z_y)}{\sqrt{E(Z_y)}}$$

$$z(y) = \frac{f(y) - E(Z_y)}{\sqrt{E(Z_y)(1 - \hat{p})}}$$

are monotonically **increasing** with $|y|$,
for all y in class C

Theorem

If $p_{\max} < \min \left\{ 1/\sqrt{4|y|}, \sqrt{2} - 1 \right\}$, then

$$z(y) = \frac{f(y) - E(Z_y)}{\sqrt{\text{Var}(Z_y)}}$$

is monotonically **increasing** with $|y|$,
for all y in class C

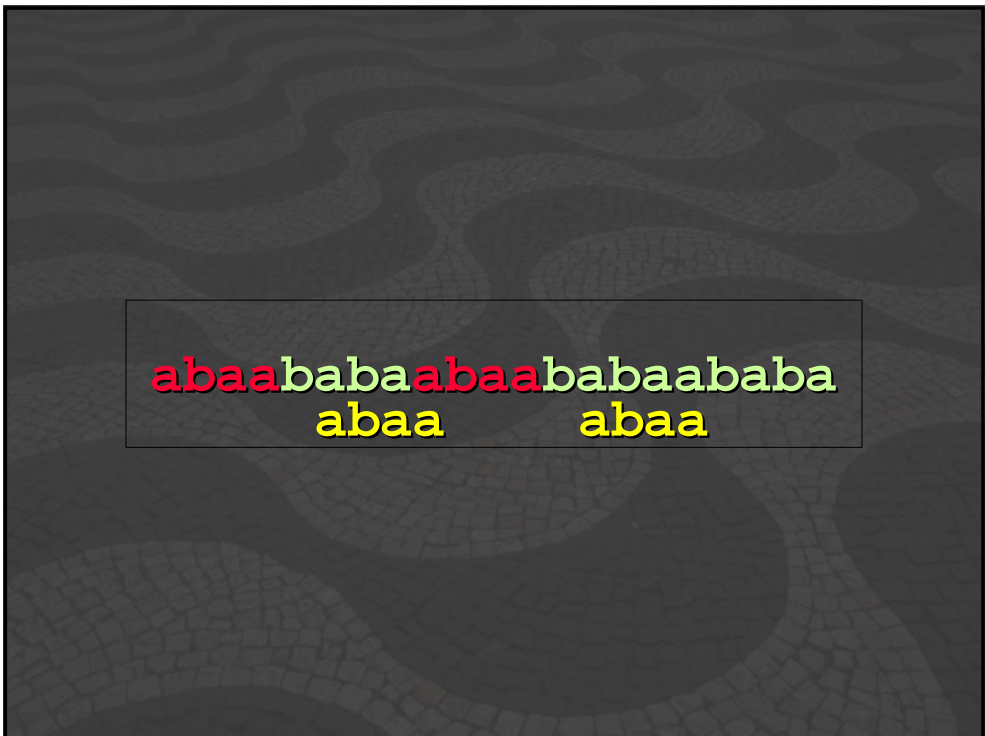
Building the partition

abaababaabaababaababa

aba**aa**babaaba**aa**babaababa
aa aa



a**baa**baba**baa**babaababa
baa baa

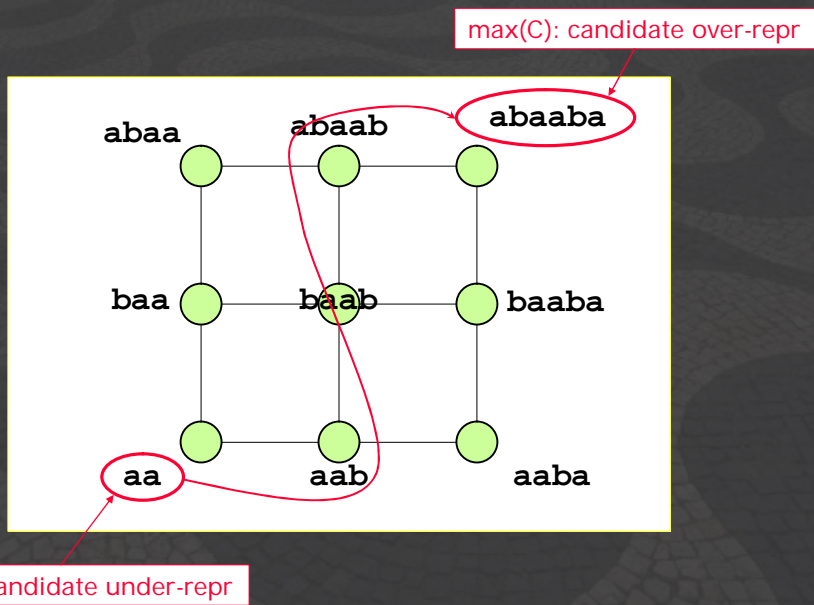


abaababa**abaa**babaababa
abaa abaa

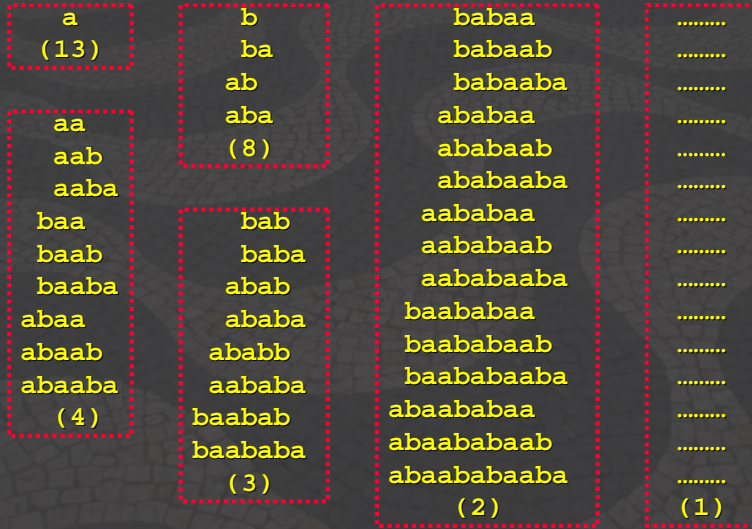
abaababaabaabaabaaba
abaab abaab

abaabaabaabaabaaba
abaaba abaaba

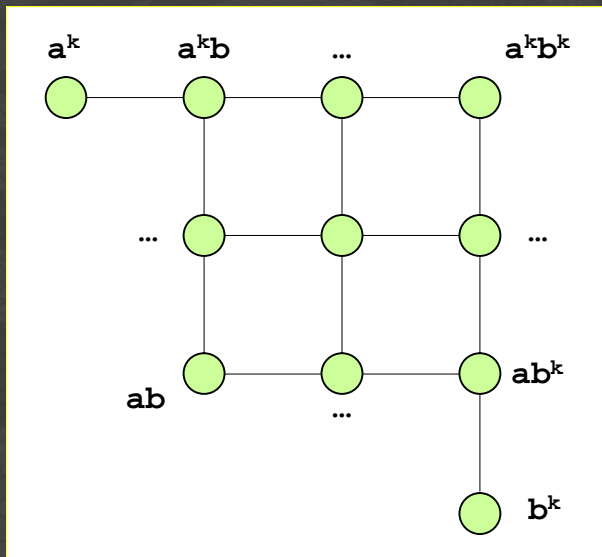
abaaba baaba ababa ababa
abaaba abaaba



$$x = \text{abaababaabaababaababa}$$



$$x = a^k b^k$$



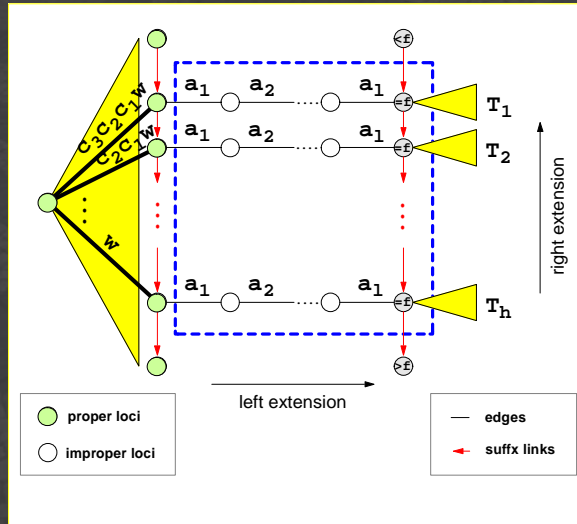
The partition $\{C_1, C_2, \dots, C_l\}$ of the set of all substrings of x , has to satisfy the following properties

- $\min(C_i)$ and $\max(C_i)$ are unique
- all w in C_i belong to some $(\min(C_i), \max(C_i))$ -path
- all w in C_i have the same count for all $1 \leq i \leq l$.

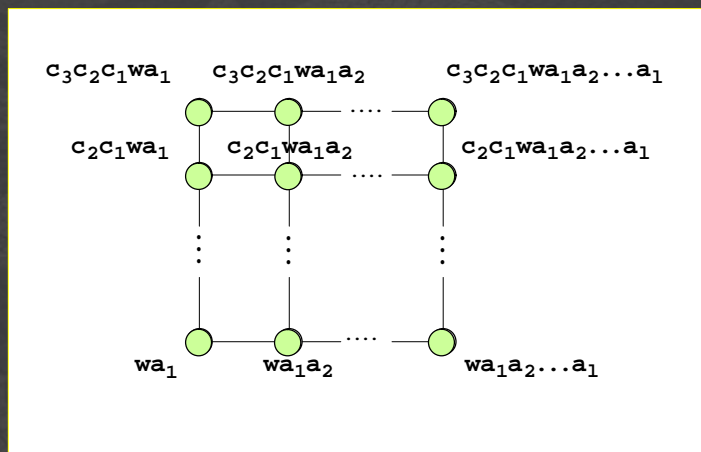
Suffix trees

- Suffix trees can be built in $O(n)$ time and space [W73,M76,U95,F97]
- Number of occurrences can be computed in $O(n)$ time

Finding equivalence classes



Finding equivalence classes



Number of classes

Theorem

The number of classes is at most $2n$

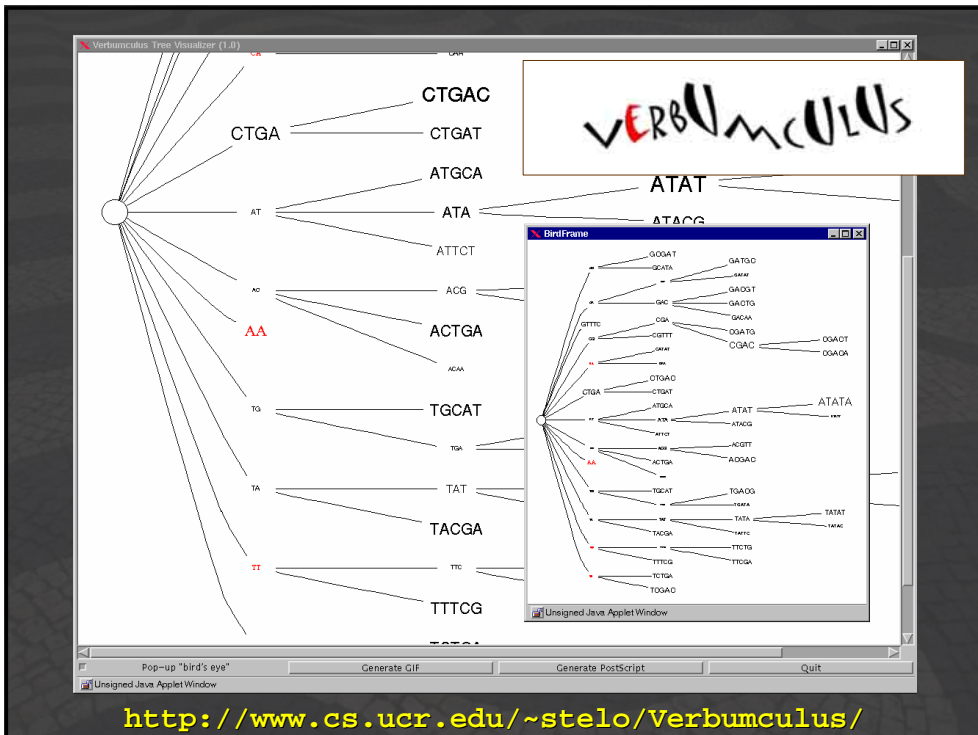
Algorithm

- Find the $O(n)$ equivalence classes
- Compute expectation, variance and score on two words (*candidates*) in each equivalence class
- Visualize the scores of the candidates

Overall time/space complexity

Theorem:

The set of over- and under-represented words can be detected in $O(n)$ time and space

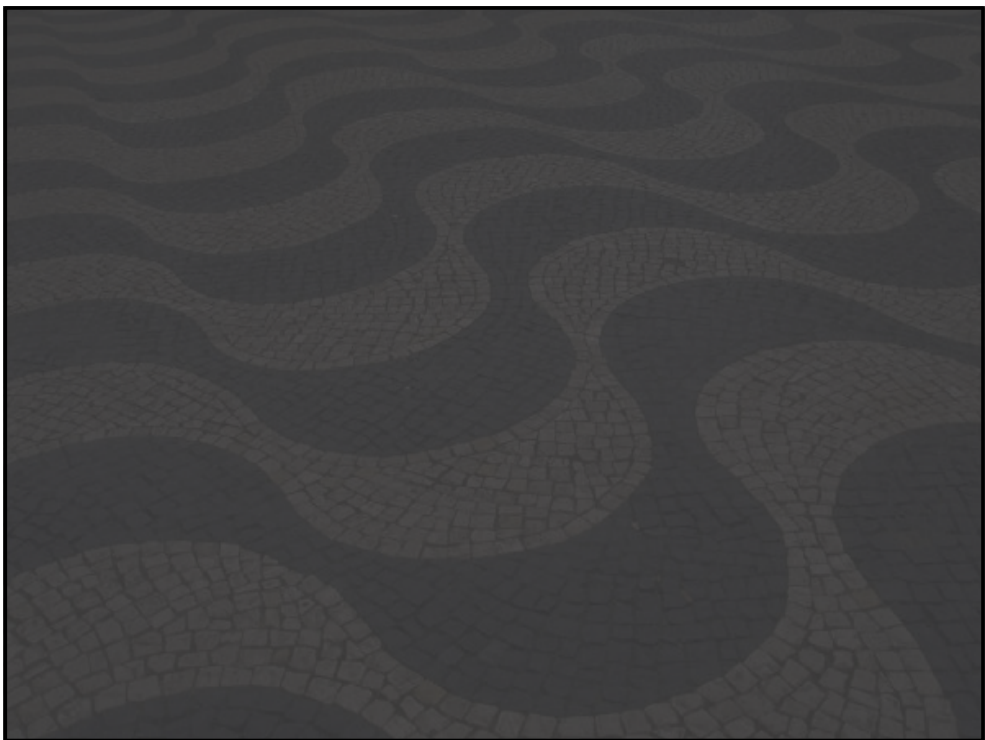
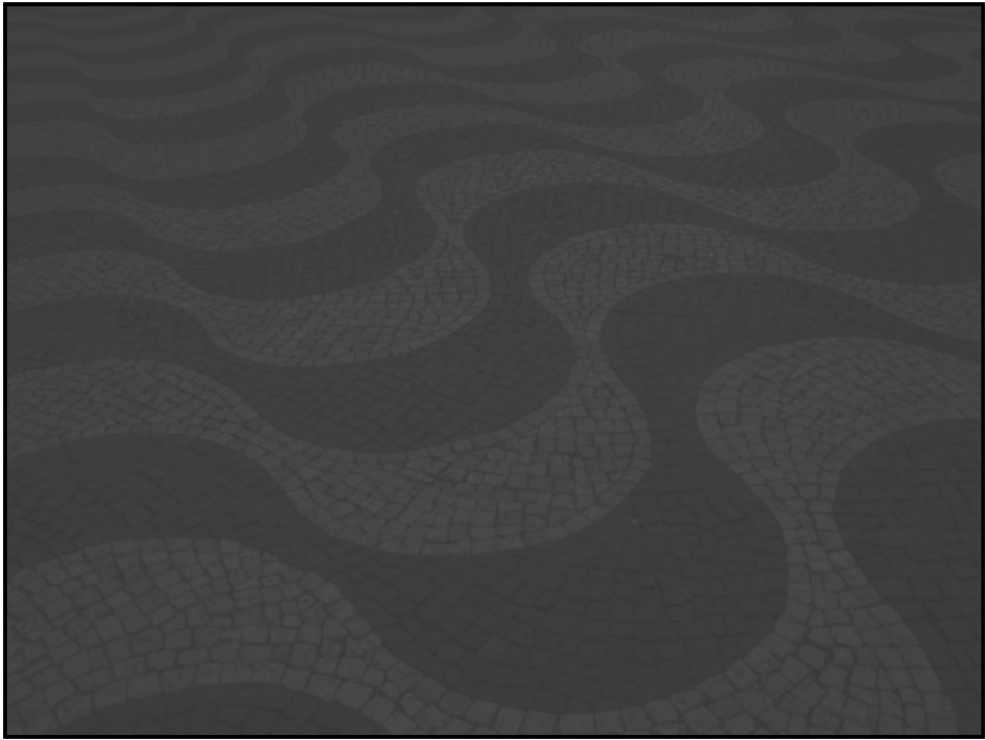


Conclusions

- Counts, expectations, variances and scores can be computed in *linear* time
- Exact patterns can be “discovered” in *linear* time and space
- Markov models and other types of counts can be handled within the same time-complexity

References

- “Monotony of Surprise and Large-Scale Quest for Unusual Words”, *RECOMB*, 2002, with A.Apostolico and M.E.Bock (to appear)
- “A Speed-up for the Commute between Subword Trees and DAWGs”, *Information Processing Letters*, 2001, with A.Apostolico (to appear)
- “Efficient Detection of Unusual Words”, *Journal of Computational Biology*, vol.7(1/2), 2000, with A.Apostolico, M.E.Bock and X.Xu
- “Linear Global Detectors of Redundant and Rare Substrings”, *IEEE Data Compression Conference*, 1999, with A.Apostolico and M.E.Bock



Text "events"

- Occurrences
 - distance constraints (non-overlapping, adjacent, max distance, ...)
 - sliding window
 - ...
- Colors
- ...

Exact or approximate?

Bernoulli Model (colors)

Let W_y be a r.v. for the number of colors of y in $\{x_1, x_2, \dots, x_k\}$, and $E(Z_y^i)$ be the expected number of occurrences of y in the i -th sequence ($1 \leq i \leq k$),

$$\bullet E(W_y) = \sum_{i=1}^k \left(1 - e^{-E(Z_y^i)} \right)$$

Scores based on colors

$$z_7(y) = c(y) - E(W_y)$$

$$z_8(y) = \frac{c(y) - E(W_y)}{\sqrt{E(W_y)}}$$

$$z_9(y) = \frac{(c(y) - E(W_y))^2}{E(W_y)}$$

where W_y is a r.v. for the number of colors of y

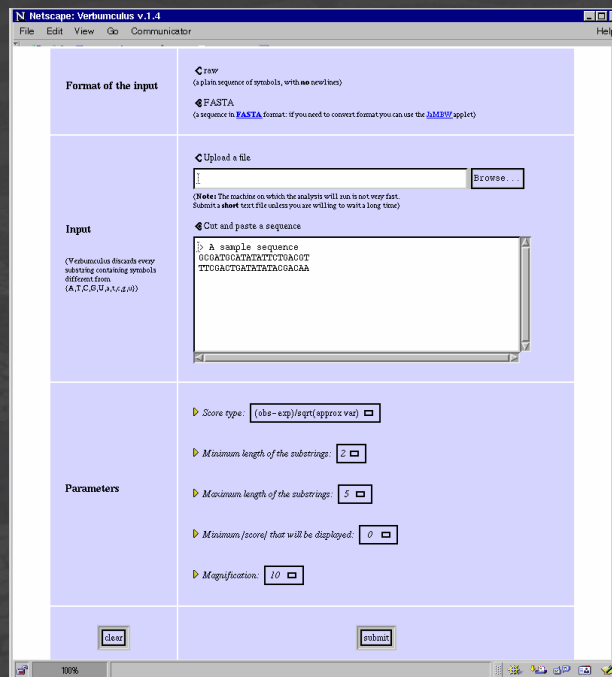
Main result

An efficient algorithm for the problem of detecting words that are, by some statistical measure, surprisingly frequent or rare in the context of larger sequences

The logo for Verbumculus, featuring the word in a stylized, hand-drawn font. The 'V' is red, and the 'U's are black.

<http://www.cs.ucr.edu/~stelo/Verbumculus/>

- Verbumculus = Verbum + Dot + TreeViz
- Verbum builds and annotates the tree
- Dot and TreeViz draw the tree; the font size of the labels is **PROPORTIONAL** to the score
- C++/STL + Perl + Java \approx 15,000 lines
- Solaris/Linux

A screenshot of the Verbumculus v.1.4 web interface running in a Netscape browser. The interface is divided into several sections: 'Format of the input' with radio buttons for 'raw' and 'FASTA'; 'Input' with a file upload button and a text area containing a sample sequence; and 'Parameters' with input fields for 'Score type', 'Maximum length of the substrings' (set to 2), 'Maximum length of the substrings' (set to 5), 'Maximum /score/ that will be displayed' (set to 0), and 'Magnification' (set to 10). There are 'Done' and 'Submit' buttons at the bottom.

Hypothesis: “Unusually frequent” patterns in the upstream sequence of a set of *co-expressed* genes are plausible binding sites implicated in transcriptional regulation

Sets of *co-expressed* genes can be identified, e.g., by DNA microarray experiments

Pattern Discovery Tools

- Exact patterns: Yeast-Tools, R'MES, WordUp (GCG), ...
- Flexible patterns: MEME (UCSD), YEBIS, SPEXS (EBI), Gibbs Sampler, BlockMaker, Teiresias (IBM), PRATT, Consensus, Winnower (UCSD), Projection (UW), ...

Typical algorithms

Naïve approach ☆

Enumerate and test all words composed by l symbols, for $1 \leq l \leq n$

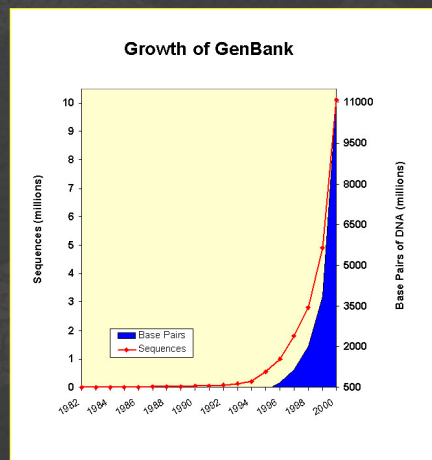
Naïve approach 🕒

Enumerate and test all words which occur in the sequences

Biomolecular Databases

- Massive
- Growing exponentially

Example: GenBank contains approximately 11,720,000,000 bases in 10,897,000 sequence records as of February 2001



$n = 1,000,000$ $S = \{A, C, G, T\}$

Naïve approach ☆

Words to be tested $O(|S|^n)$

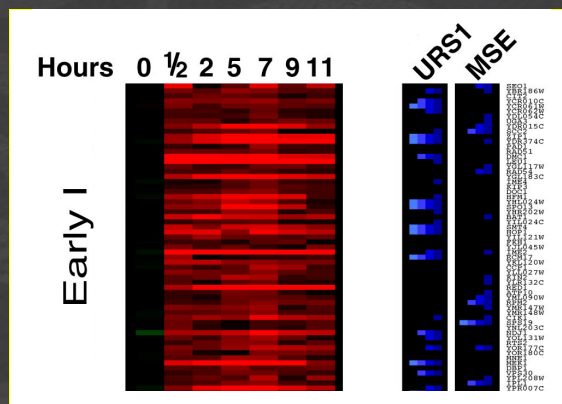
in this case $\propto 4^{1,000,000}$

Naïve approach 🕒

Words to be tested $O(n^2)$

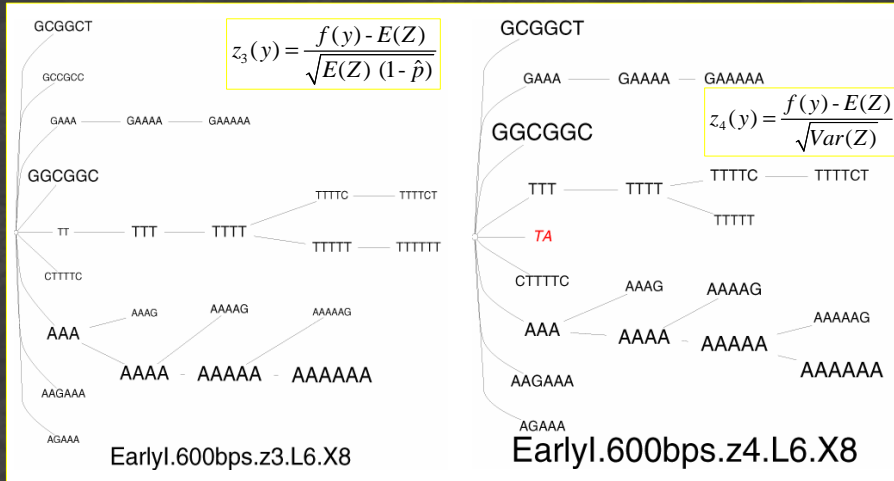
in this case $\propto 1,000,000^2$

Cluster Early I

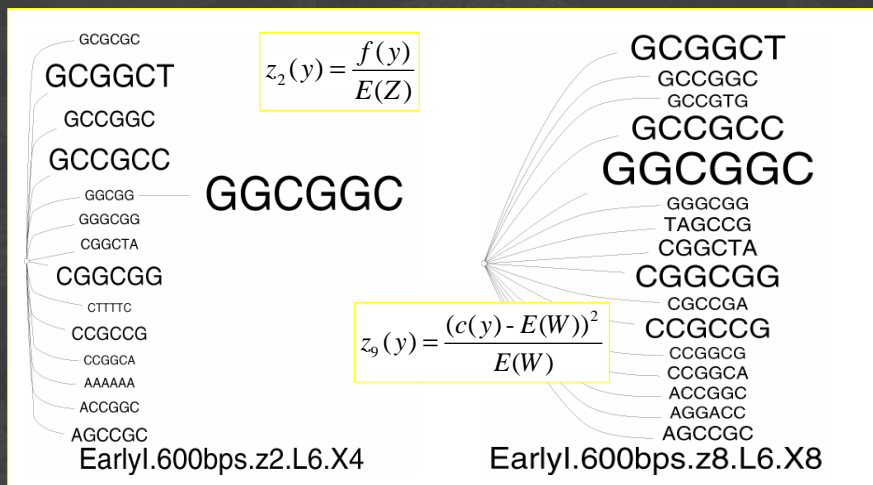


Dataset from "The Transcriptional Program of Sporulation in Budding Yeast", by S.Chu, J.L.DeRisi, M.B.Eisen, J.Mulholland, D.Bodstein, P.O.Brown, I.Herskowitz, *Science*, 1998

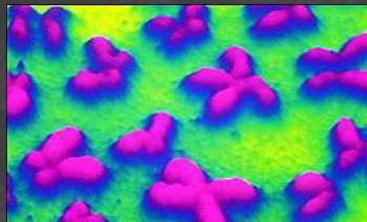
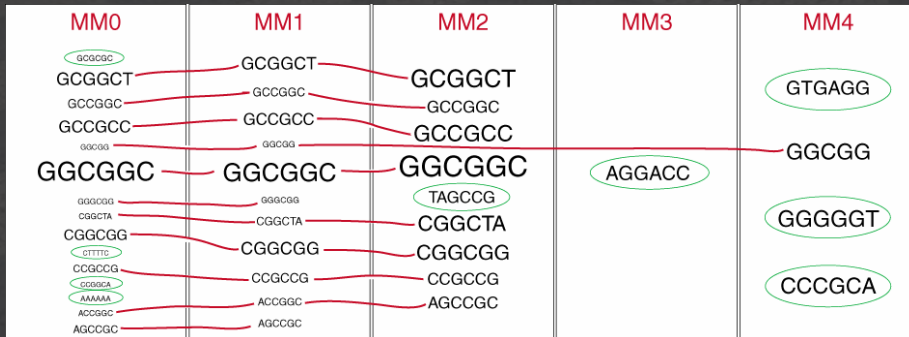
Analysis of EarlyI (1/3)



Analysis of EarlyI (2/3)



Analysis of EarlyI (3/3)



Organism: E. coli K12
number of strands = 2025
number of bases = 1992558
number of 4-grams checked (overlapping) = 1787476
expected frequency (uniform distribution) = 6982.33

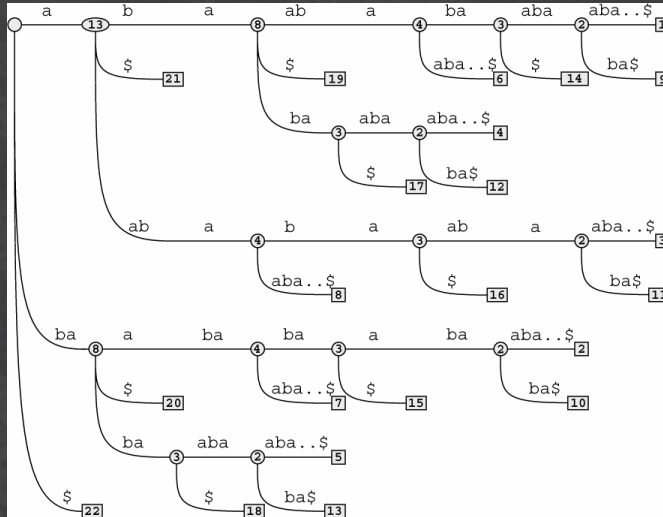
4-gram	f(y)	f(y)/total	f(y)/exp
CTAG	99	0.0001281136	0.0327970837
TAGG	99	0.0005577697	0.1427890500
ATAG	1262	0.0007060235	0.1807420072
TAGA	1272	0.0071161799	0.1821741942
TAGT	1361	0.0007614088	0.1949205591
CCTA	1605	0.0008791422	0.2298660234
CCCC	1660	0.0009286718	0.2377430522
GAGG	2055	0.0011496658	0.2943144411
TTAG	2199	0.0012302263	0.3149379348
CTAT	2337	0.0013074301	0.3347021163
TATG	2372	0.0013270108	0.339717710
TATA	2433	0.0013611372	0.3484511131
CTAA	2461	0.0013768017	0.3524612258
TAGC	2574	0.0014400193	0.3696449496
GTAG	2609	0.0014596000	0.3736576044
TCTA	2658	0.0014870130	0.3806753210
GCCC	2801	0.0015670140	0.4011555959
CCTC	2833	0.0015849164	0.4057385947
AGAC	2970	0.0016615608	0.4253591713
ACTA	3007	0.0016822603	0.430786494
AGTC	3144	0.0017589047	0.4502796121
CCCA	3154	0.0017644992	0.4517117992
AGTA	3208	0.0017947094	0.4594456093
CTCC	3236	0.0018107748	0.4634557311
AGCG	3278	0.0018318708	0.4694709188
TCCC	3282	0.0018361086	0.4700437936
TGTA	3326	0.0018607243	0.4763454167
CCTC	3350	0.0018741510	0.4797826656
GAGT	3407	0.0019032423	0.4872300383
GGA G	3426	0.0019166691	0.4906672873
CTTA	3429	0.0019183474	0.4910969434
CTTG	3454	0.0019323336	0.4946774111
CAAG	3493	0.0019541521	0.5002639406
ATTA	3543	0.0019821245	0.5074238759
GATA	3553	0.0019877190	0.5088560630
CGAG	3554	0.0019882784	0.5089992017
AGGA	3559	0.0019910757	0.5097153755
ACTC	3657	0.0020459016	0.5237501734
AGAG	3692	0.0020654823	0.528734631
CTCA	3755	0.0021007275	0.537862416
TAA T	3756	0.0021012870	0.5379294603
CACA	3780	0.0021147137	0.5413667093
GACAC	3924	0.0021955712	0.5619902029
CCTT	3932	0.0021971498	0.5631359526
GCGG	3935	0.0022014282	0.5635656087
ACAC	3988	0.0022107939	0.5711562001
GACT	4023	0.0022506596	0.5761688549
ACTT	4035	0.0022573730	0.5778874793
TACA	4077	0.0022808698	0.5839026650
G TGT	4111	0.0022998910	0.5887210110
G GGA	4156	0.0023250662	0.5952169428
CCTT	4229	0.0023659059	0.6056719083
T CCT	4246	0.0023754165	0.6081066263
TAA T	4380	0.0024503826	0.6272979330
CTCA	4380	0.0024503826	0.6272979330
GCTC	4454	0.0024917817	0.6378961172
TGAG	4493	0.0025136002	0.6434816467
TCTT	4503	0.0025191947	0.6449138338
ACTT	4510	0.0025231108	0.6459163648
G GGT	4556	0.0025488454	0.6525044252
CTAT	4580	0.0025622722	0.6559416742
GCCC	4620	0.0025846501	0.6616704224
ATAA	4698	0.0026282870	0.6728414815
TGTC	4750	0.0026575783	0.6802885542
GCTA	4751	0.0026579378	0.6804320729
CTAT	4774	0.0026900567	0.6807185103
GACA	4795	0.0026925535	0.6867336900
TCTC	4807	0.0026992669	0.6884523205
ATAA	4824	0.0027069877	0.6908870385
AGGT	4910	0.0027689900	0.7032038472
CACA	4928	0.0027567701	0.7057817839
CACT	4936	0.0027614357	0.7069275336
ACCC	4967	0.0027787786	0.7116731135
AGTT	5046	0.0028229750	0.7216815912
CCTG	5047	0.0028235344	0.7221748100
TGTT	5112	0.0028598985	0.7321347159
T CAT	5151	0.0028817170	0.7377195557

Definition:

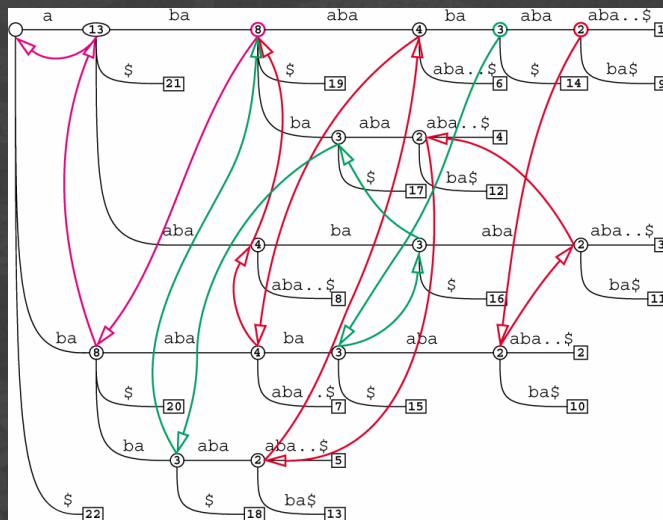
Given a substring w of x the **implication** of w in x , denoted by $imp_x(w)$, is the string $u w v$, such that

- every time w occurs in x , it is preceded by u and followed by v
- u and v are maximal

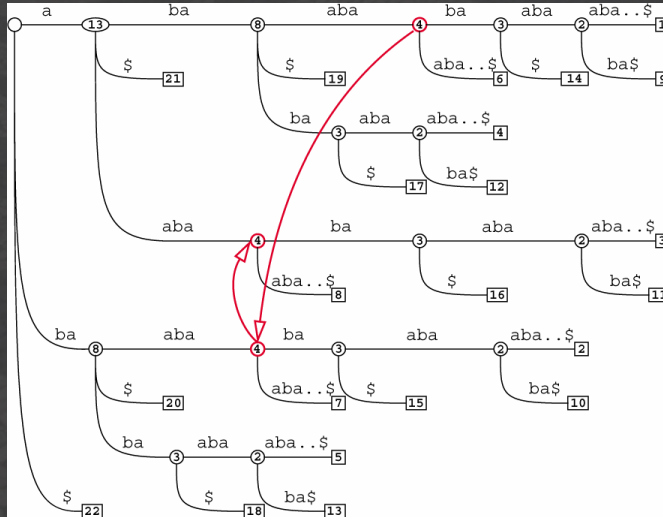
Finding Equivalence Classes



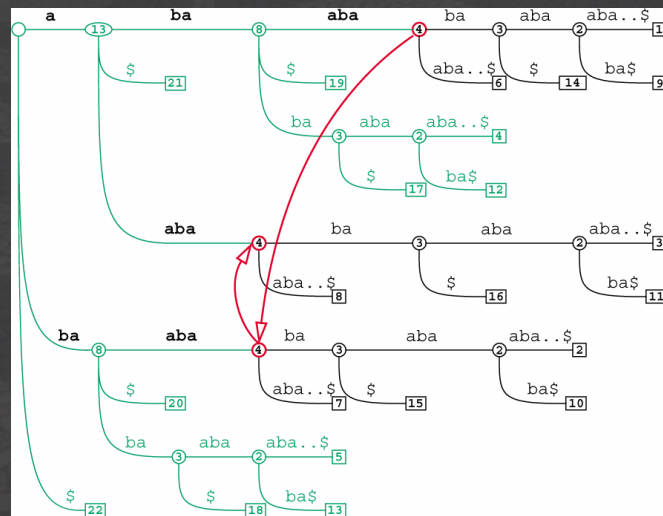
Finding Equivalence Classes



Finding Equivalence Classes



Finding Equivalence Classes



What's next?

- extension to other types of count and hidden Markov models
- estimation of statistical parameters by "shuffling"
- more experiments on biosequences and in other domains
- extension to approximate/flexible patterns

How to choose the threshold

$$P\left(\left|\frac{f(y) - E(Z_y)}{\sqrt{\text{Var}(Z_y)}}\right| > 2\right) = .0456$$

