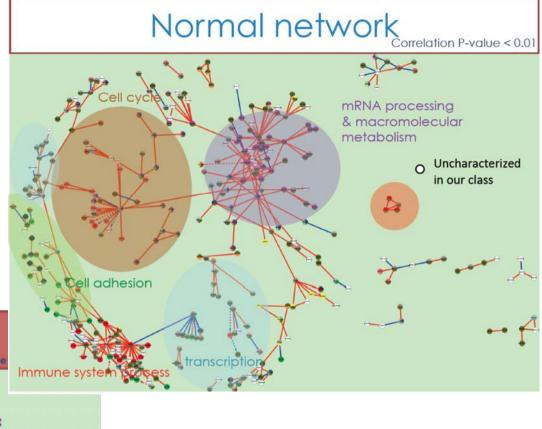
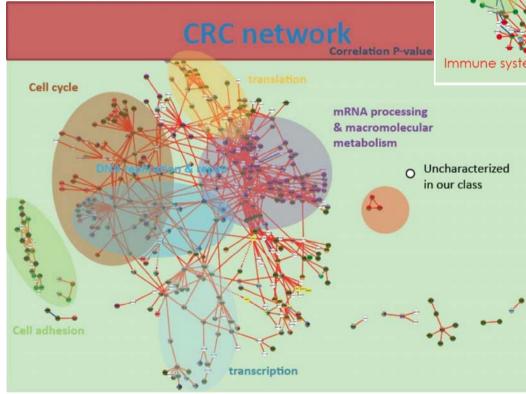


Graduate Institute of Systems Biology National Central University Zhongli/Jhongli/Chungli, Taiwan

Enriched analysis of microarrays





### Three questions I like to ask

- WHAT is the phenomenon?
  - What is important/unusual/interesting?
- HOW did it happen?
  - (Physics)
- WHY did it happen?
  - (Biology)

### What we plan to discuss

How did life evolve so rapidly • random sequence • statistical properties of the genome • uniformity homogeneity and universality order of genome o genome the information carriergenome at the "edge of chaos" • symmetries in the genome • coding versus non-coding regions • how genome grew • segmental & whole genome duplication • genome the blind self-plagiarizer
 genome in a state of fixed-point • the nearly universal cumulative point mutation density • age of the genome • the mystery of the "missing" mutations • • the phyletic gradualism versus punctuated equilibria debate

**Motivation**: Understand large scale growth and evolution of genome; how did it evolve so rapidly

Materials: Complete genome of prokaryotes and eukaryote

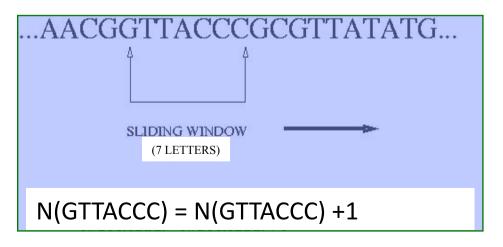
**Method**: Variety of statistical analysis of frequency of oligonucleotides (k-mers)

### Sequence alignment versus statistical analysis

Feature	Seq alignment	Stats analysis	
Nature of quantity measured	relative	absolute	
Sequence length	< 10k	> 5k	
Sequence comparison relative word in order word content	excellent poor	poor moderate	
Positional information local, relative large scale structure	excellent poor	poor good	
Correlation short range long range	good poor	poor good	
(Absolute) Order of sequence		excellent	
<pre>Evolution time scale   &lt; 100 Mya (divergence time)</pre>	good 	 good	
Mode of genome growth		good	

### Methods

- Treat genome as a 4-letter pseudo-text
- Count frequency of occurrence of "words"
  - A k-letter word is called a "k-mer"
- Number of possible k-mers is  $4^k$
- Use a sliding window of width k and slide 1
- Analyze set of k-mer frequencies



After a single sweep have, for each k, a  $4^k$ -component vector  $\{f_i \mid i=4^k\}$  whose element  $f_i$  is the frequency of the ith k-mer.

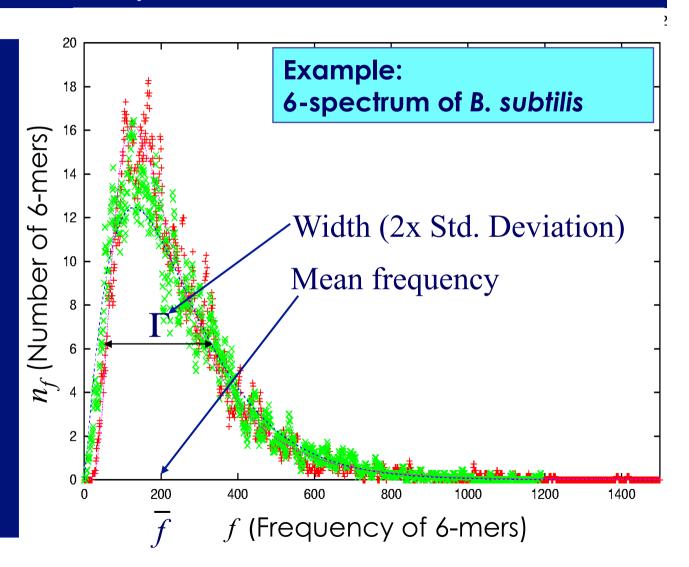
## Frequency set, *k*-spectrum & relative spectral width

Given freq. set  $\{f_i\}$ , define

k-spectrum  $\{n_f \mid f=1,2,...\}$  $\Sigma_i f_i = \Sigma_n f n_f$ 

Relative spectral width

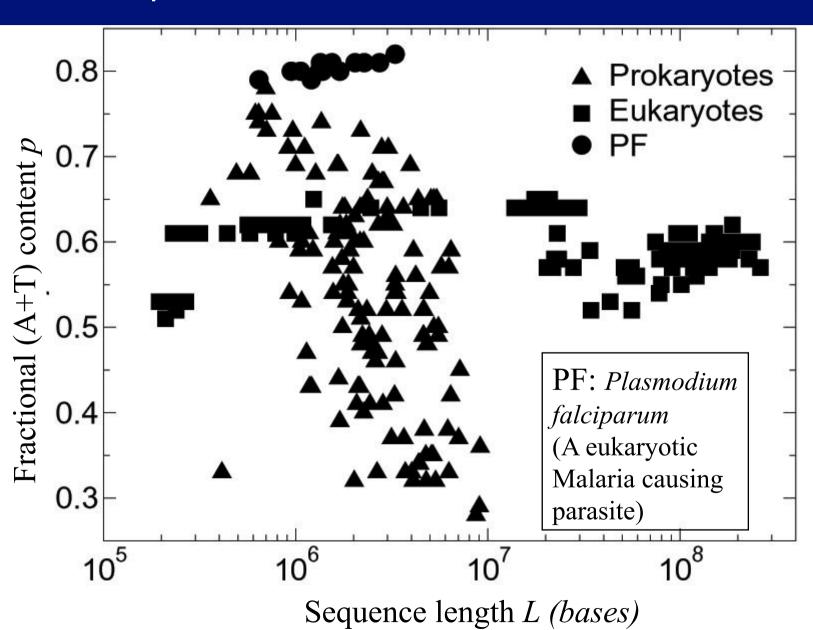
 $\sigma$  = std dev/<f>



### Data set

- All sequences downloaded from GenBank (Nov 2006)
- 422 complete prokaryotic ("bacterial") genomes
  - Lengths 0.5 to 9 M base pairs (Mb)
- 19 complete eukaryotic genomes totaling 328 complete chromosomes
  - Lengths 0.2 to 250 Mb
- Total data set: 2.2 x 10<sup>10</sup> base pairs

### Complete Genomes are diverse



### Four lectures

**Lecture One**. Genome: A Large System with Small System Statistics

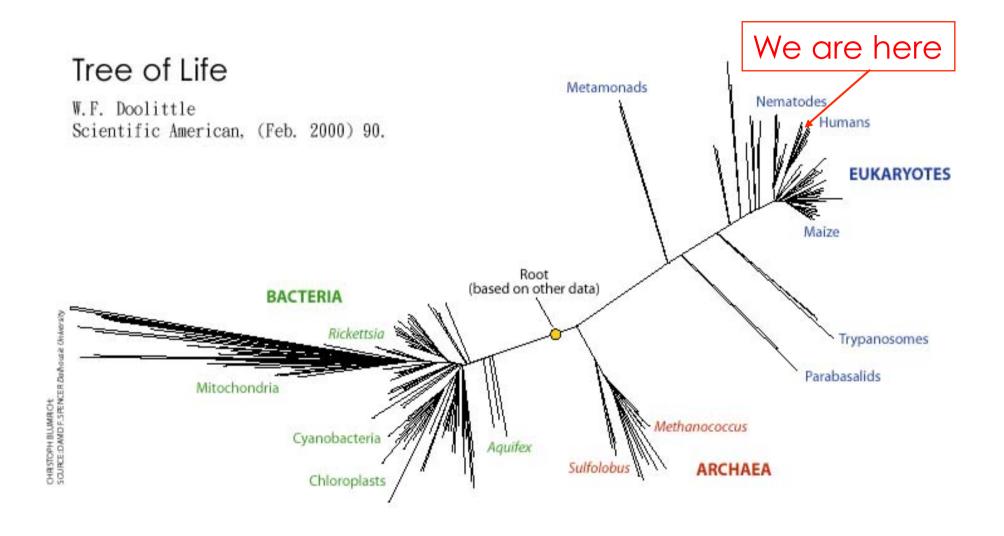
**Lecture Two**. Genome: Information Carrier at the "Edge of Chaos"

Lecture Three. Symmetries in Genomes

**Lecture Four**. Genome the Blind Self-Plagiarizer and the Mystery of the "Missing" Mutations



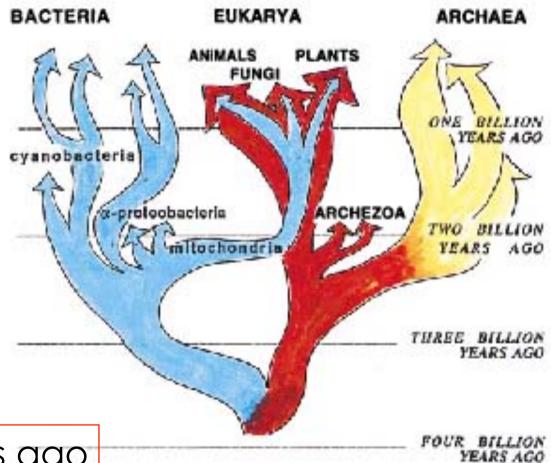
### Life is highly diverse and complex



### And it took a long time to get here

Divergence of species W.F. Doolittle, PNAS 94 (1997) 12751.

now



4 billion yrs ago

### Life evolved super rapidly

- Typical mutation rate is 1/site/Byr
- Life in the form of microbe existed less than
   1 Byr after cooling of earth
- Suppose we were "given" a genome 1000 bp long at dawn of life
- After 1 Byr every site mutated once
- Size of sequence space is  $4^{1000} \sim 10^{600}$
- Evolution could not have been driven by point mutation

## Evolution of Genomes and the Second Law of Thermodynamics

### Genomes grew & evolved stochastically

- modulated by natural selection
- Bigger genomes carry more information than smaller ones

### The second law of thermodynamics:

- the entropy of closed system can never decrease
- a system that grows stochastically tends to acquire entropy
- Increased randomness more entropy

#### Shannon information

Information decreases with increasing entropy

### Long-range variation in GC content

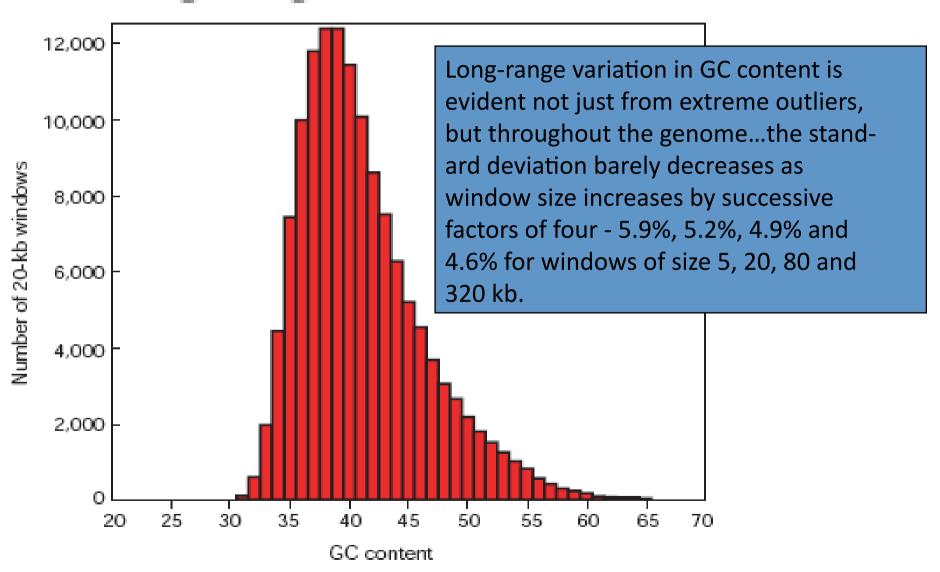
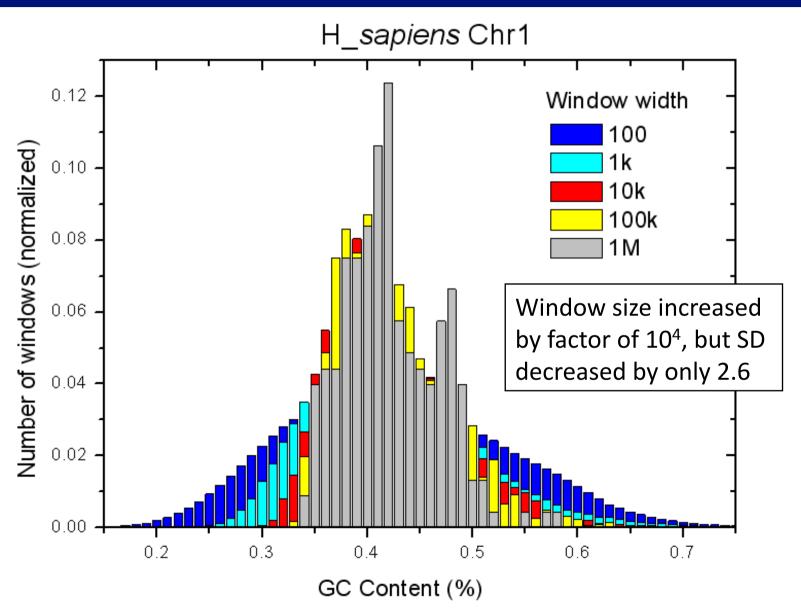
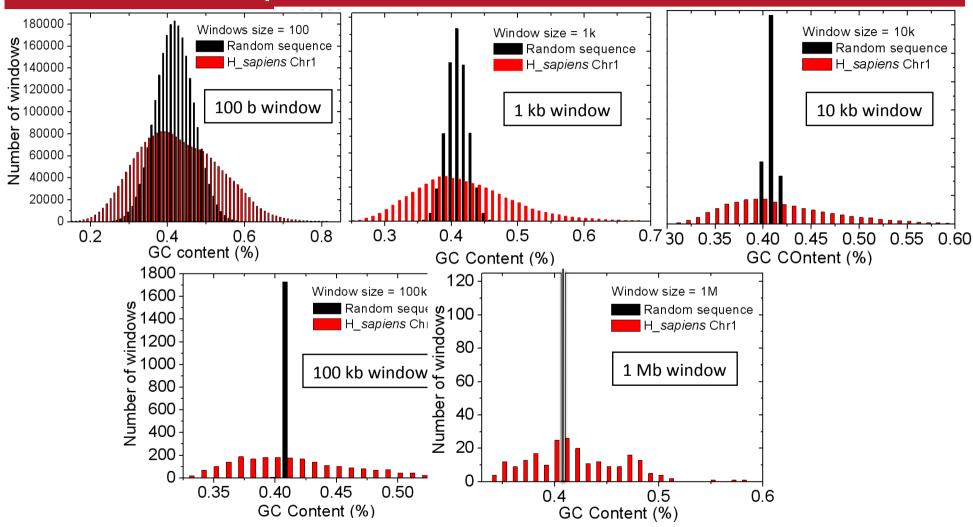


Figure 12 Histogram of GC content of 20-kb windows in the draft genome sequence.

## CG-content in Human genome has long-range variation

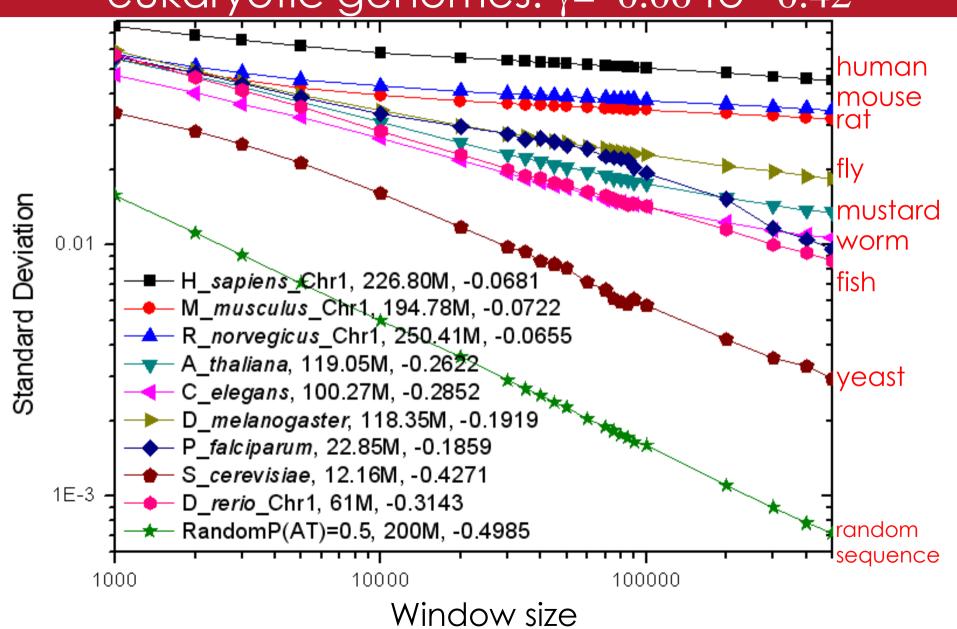


# For GC-content histograms: sample size = window size



Variation of CG-content in Human genome does not obey central limit theorem

### Power law is universal for complete eukaryotic genomes: $\gamma$ =-0.06 to -0.42



#### Distribution and Width

- Consider  $\tau$  equally probable events occurring a total of L times.
- Distribution of occurrence frequency characterized by
  - mean frequency:  $f_{ave} = L/\tau$
  - SD (standard deviation)  $\Delta$ ; or CV (coefficient of variation) =  $\Delta/f_{ave}$
  - Higher moments of distribution

### Random events

- Random events given by Poisson distribution
  - $-\Delta^2 = f_{ave'}$  or,  $(CV)^2 = 1/f_{ave}$
  - That is,  $(CV)^2 = \tau/L$
- For fixed  $\tau$ ,  $(CV)^2 \sim 1/L$ 
  - Large L limit (thermodynamic limit): L ~ infinity, CV ~ 0
- For given  $\tau$ , if CV is known, then
  - $L \sim \tau/(CV)^2$

## Frequency set, "k-spectrum" & relative spectral width

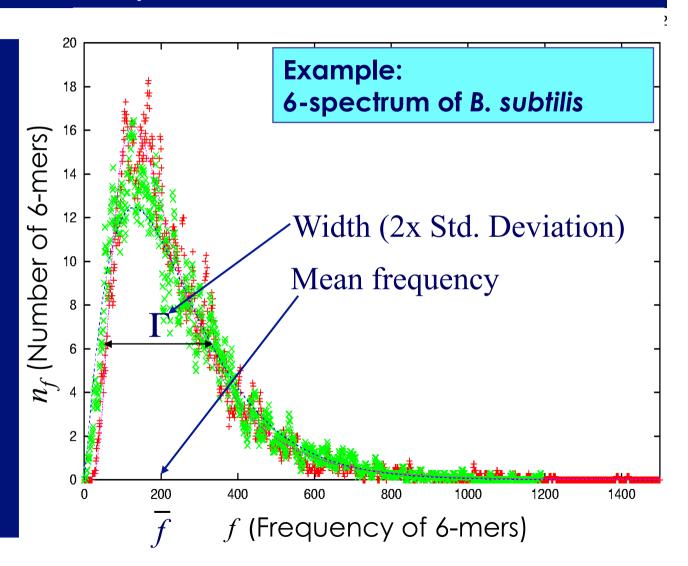
Given freq. set  $\{f_i\}$ , define

k-spectrum  $\{n_f | f=1,2,...\}$  $\Sigma_i f_i = \Sigma_n f n_f$ 

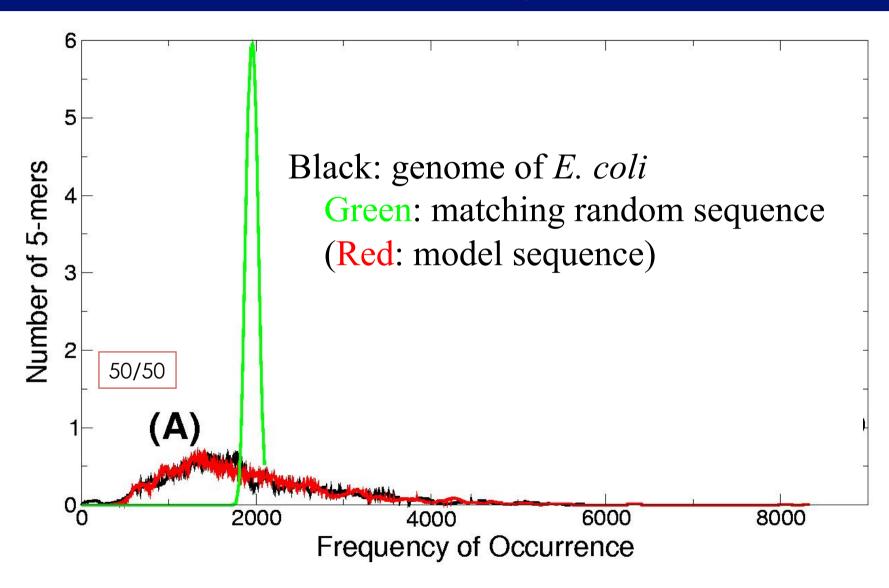
----

Relative spectral width

 $\sigma$  = std dev/<f>



# Huge difference between genomes and random sequences

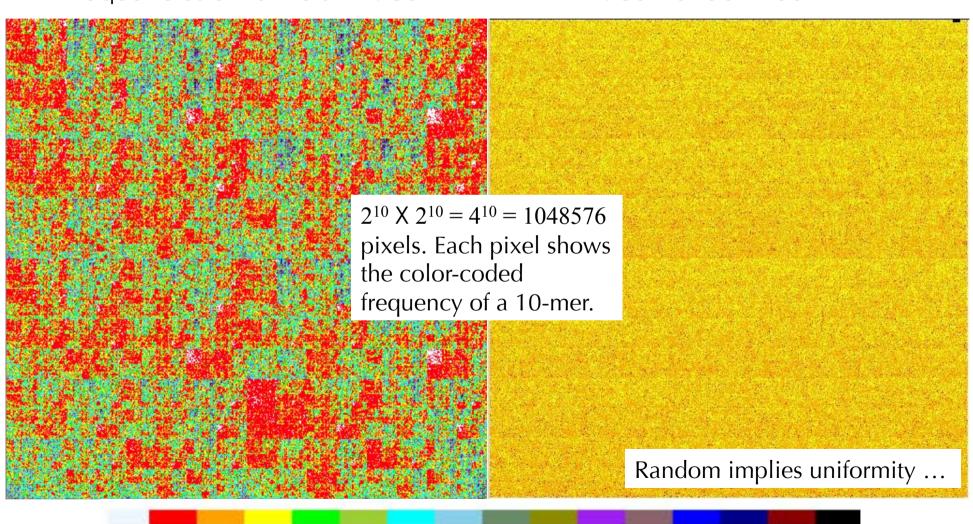


### 2D "portrait" of E. coli genome

Frequencies of 10-mers in E. coli

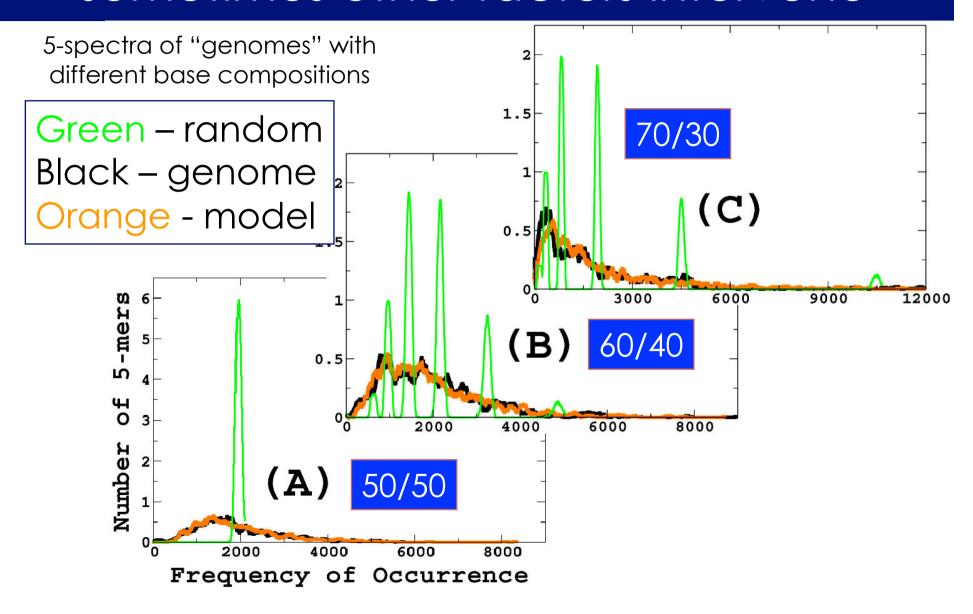
E. coli randomized

24

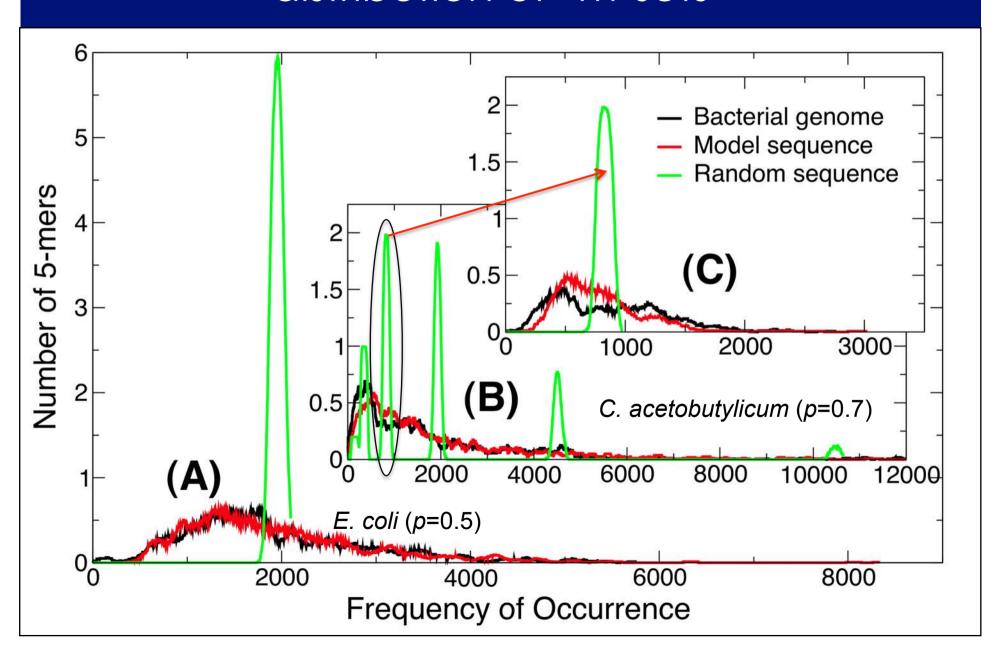


12

## Random implies uniformity; but sometimes other factors intervene



### k-spectrum is composed of (k+1) subdistribution of "m-sets"



#### Partition of k-mers into m-sets

- Consider genome with fractional AT-content p (then fractional GC-content q=1-p
- There will be more AT-rich words than GC-rich words
- Partition k-mers into sets, called m-sets, m=0,1,...,k; k-mers in an m-set all have m ATs.
- Total number of:
  - kinds of k-mers:  $\tau = 4^k$
  - k-mers: L (sequence length); mean frequency is  $L/\tau$
  - kinds of k-mers in m-set:  $\tau_m$
  - k-mers in an m-set:  $L_m$ ; mean frequency is  $f_m = L_m/\tau_m$

$$\tau_m = \binom{k}{m} 2^k, \quad L_m^{\{\infty\}} = L \binom{k}{m} p^m q^{k-m},$$

$$\bar{f} = L/\tau, \quad \tau = 4^k \qquad \bar{f}_m^{\{\infty\}} = L_m^{\{\infty\}}/\tau_m$$

### $f_m$ well approximated by its large-L limit

( <i>k</i> =5)	fm					
Sequence	(m=) 0	1	2	3	4	5
p = 0.492						
E coli	2509	2245	1877	1760	1944	2656
Random	2101	2044	1987	1922	1857	1795
$\lim_{L  o \infty} Random^*$	2114	2048	1983	1920	1860	1801
p = 0.691						
C acetobutylicum	154	397	918	1951	4272	10300
Random	176	394	882	1970	4400	9832
$\lim_{L\to\infty}$ Random*	176	393	880	1968	4402	9845

All sequences normalized to a length of 2 Mb;  $\bar{f} = 2 \times 10^6/4^5 = 1953$ . Random means matching random sequence, or sequence obtained by scrambling the genome. \*Values of  $\bar{f}_m^{\{\infty\}}$  given by Eq. (6). doi:10.1371/journal.pone.0009844.t006

### Decomposition of standard deviation into "non-fluctuating" and "fluctuating" parts

$$\sigma^{2} = \tau^{-1} \sum_{u \in S} (f_{u} - \bar{f})^{2} = \tau^{-1} \sum_{m=0}^{k} \sum_{u \in S_{m}} (f_{u} - \bar{f}_{m} + \bar{f}_{m} - \bar{f})^{2}$$

$$= \tau^{-1} \sum_{m=0}^{k} \left( \tau_{m} (\bar{f}_{m} - \bar{f})^{2} + 2(\bar{f}_{m} - \bar{f}) \sum_{u \in S_{m}} (f_{u} - \bar{f}_{m}) + \sum_{u \in S_{m}} (f_{u} - \bar{f}_{m})^{2} \right).$$
(7)

$$\sigma^2 \equiv \sigma_{nf}^2 + \sigma_{fl}^2,$$

(generalization of parallel axis theorem)

non-fluctuating part

$$\sigma_{nf}^2 = \sum_{m=0}^k \frac{\tau_m}{\tau} (\bar{f}_m - \bar{f})^2$$
, (depends only on averages)

fluctuating part 
$$\sigma_{fl}^2 = \sum_{m=0}^k \sum_{u \in S_m} \frac{(f_u - \bar{f}_m)^2}{\tau} \equiv \sum_{m=0}^k \frac{\tau_m}{\tau} \sigma_{m,fl}^2.$$
 (10)

### Large-L limit of CV

Large-L limit for non-fluctuating part is known (for all sequences)

$$(CV^{\{\infty\}})^2 \equiv \lim_{L \to \infty} CV_{nf}^2 = \sum_{m=0}^k 2^{-k} {k \choose m} (2^k p^m q^{k-m} - 1)^2$$

$$= 2^k (p^2 + q^2)^k - 1,$$
(12)

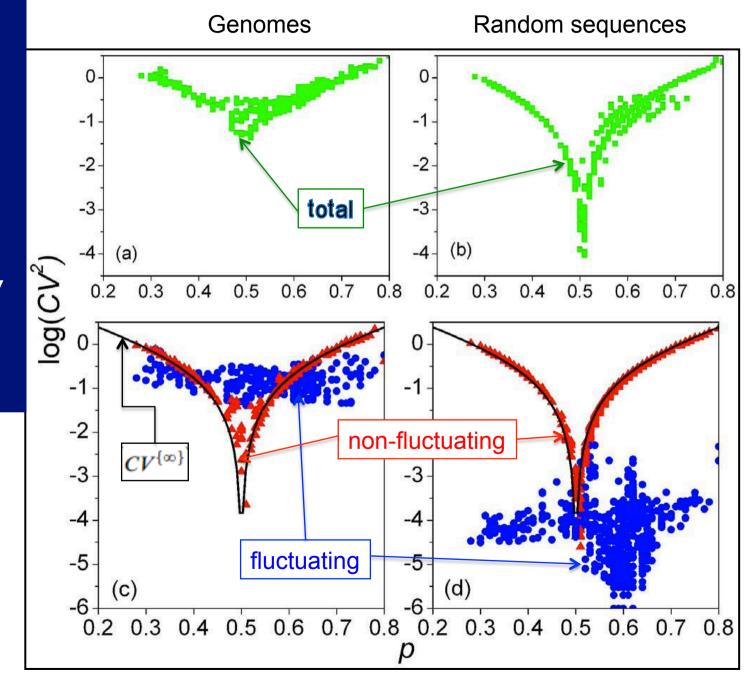
So is large-L limit for fluctuating part for random sequences

$$\lim_{L \to \infty} CV_{fl}^2 = \frac{1}{\bar{f}^2} \lim_{L \to \infty} \sigma_{fl}^2 = \frac{1}{\bar{f}^2} \sum_m \frac{\tau_m}{\tau} \bar{f}_m$$

$$= \frac{1}{\bar{f}} = \frac{\tau}{L} \qquad \text{(random sequence)},$$
(13)

For genome-size random sequence,  $CV_{nf} >> CV_{fl}$ 

Genome and random sequ's differ only in  $CV_{fl}$ 



### Equivalent length of a sequence

Because  $CV_{fl} \sim \tau/L$ , we define for any sequence an equivalent length  $I_e$ 

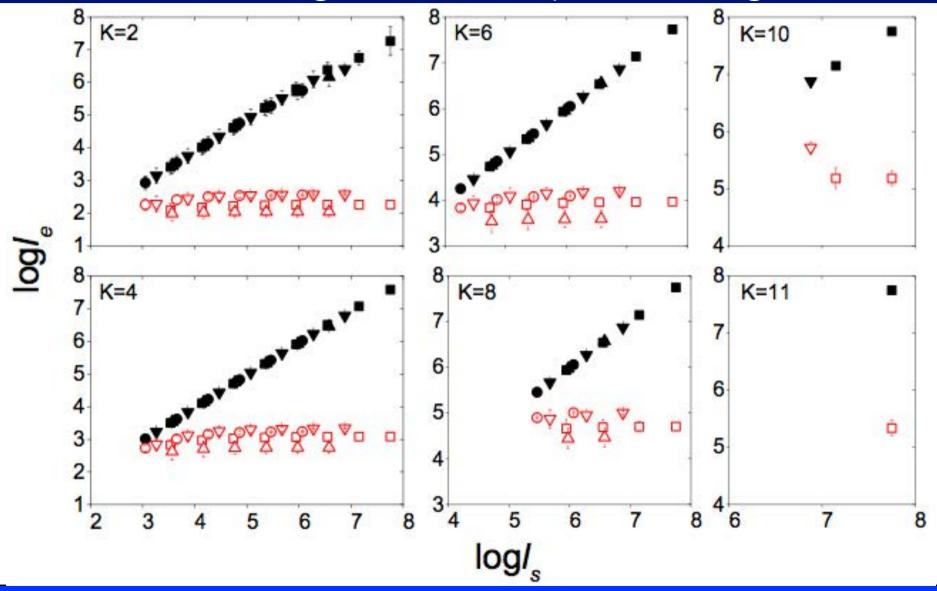
$$l_e = b_k \tau / CV_{fl}^2 \qquad \text{(for } k \ge 2\text{)} \tag{14}$$

 $[b_k = 1 - 1/2^{k-1}; = 1/2 \text{ for k} = 2 \text{ and approaches 1 quickly for larger k's}]$ 

(note: for random sequence  $L = I_e$ ). The  $I_e$  of a sequence is the length of the random sequence whose  $CV_{fl}$  is the same as that of the sequence.

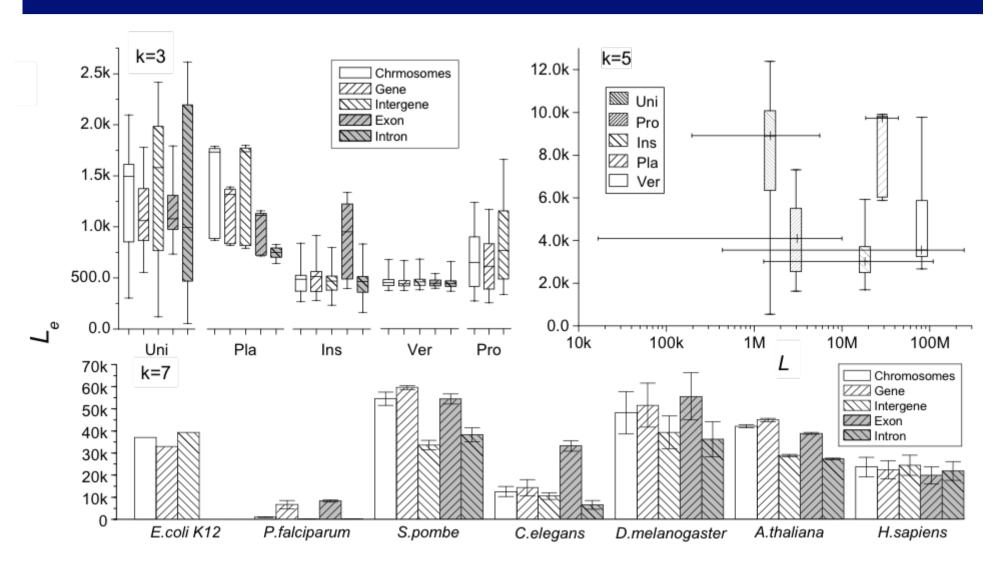
Notation: use  $l_e$  for a segment, and  $L_e$  for a whole genome/chromosome

### Genomic (E. coli, worm, mustard and human) $l_e$ does not grow with sequence length



Red: segments from genomes. Black: segments from random sequences.

## L<sub>e</sub> of coding and non-coding parts not much different (I)

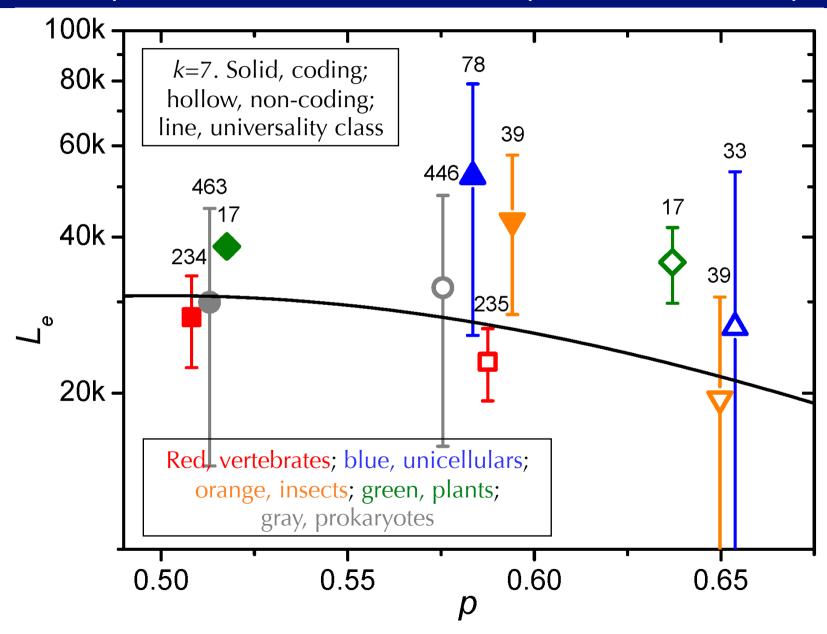


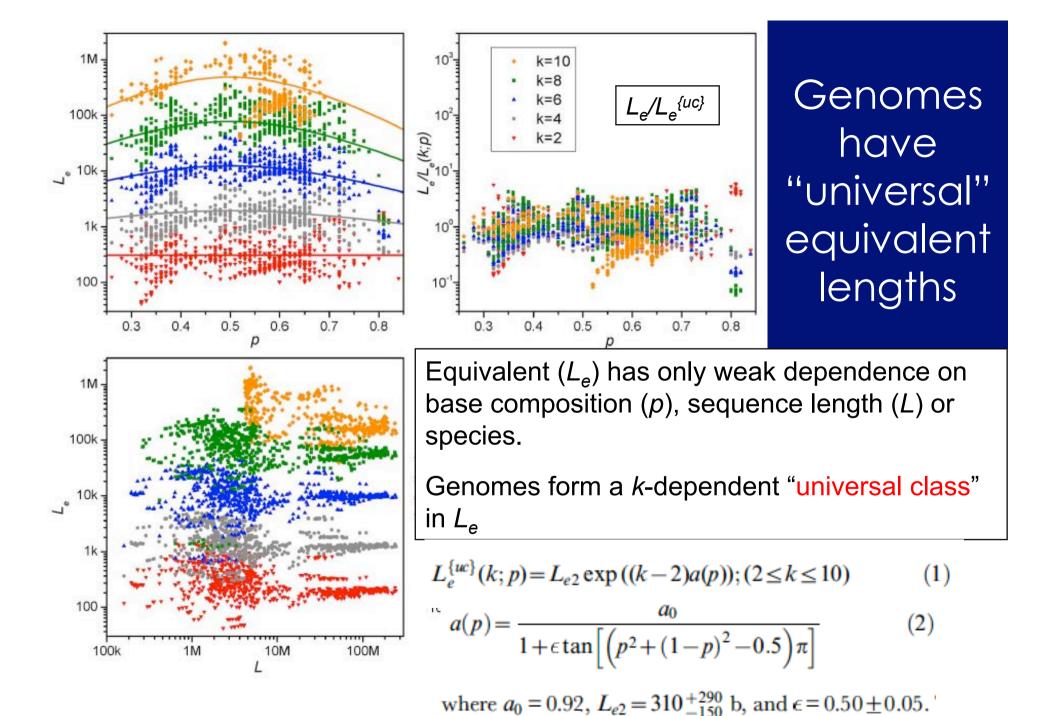
### $L_{\rm e}$ of coding and non-coding parts not much different (II)

Category	$L_e$ (kb)					
	(k =) 2	5	7	10		
All	.359+.333	$4.56^{+3.60}_{-2.01}$	33.7+30.0	388 <sup>+524</sup> <sub>-223</sub>		
gn (41.8%)	. 317+.253	$4.21^{+2.82}_{-1.67}$	$31.2^{+23.7}_{-13.4}$	337 <sup>+396</sup> <sub>-186</sub>		
ig (59.6%)	$.462^{+.879}_{302}$	$4.99^{+4.49}_{-2.36}$	$31.6^{+26.9}_{-14.5}$	$213^{+170}_{-95}$		
ex (3.3%)	.292+.215	$4.40^{+2.55}_{-1.62}$	35.3 <sup>+20.8</sup> -13.1	620+298		
in (31.8%)	$.348^{+.679}_{230}$	$3.65^{+2.55}_{-1.50}$	$23.5^{+13.9}_{-8.7}$	$213^{+206}_{-105}$		
$L_e^{\{uc\}}\ (p=0.5)$	.310 <sup>+.290</sup> <sub>150</sub>	$4.90^{+4.58}_{-2.24}$	$30.1^{+28.1}_{-13.8}$	487+455		
RSD model	.597 + .756	$4.79^{+0.82}_{-0.70}$	$32.0^{+7.0}_{-5.8}$	$510 + \frac{211}{-149}$		

 $L_e(k)$ , k=2, 5, 7 and 10, averaged over 865 chromosomes. Total sequences length is about  $2.2 \times 10^{10}$  bases. Abbreviations: All, complete chromosome; gn, genes; ig, intergenic; ex, exons; in, introns. Percentage given indicates portion of complete sequence.  $L_e^{\{uc\}}$  is defined in Eq. (1) and RSD results are averaged over 200 model sequences. See Table S4 for  $L_e(k)$  of other k values. doi:10.1371/journal.pone.0009844.t002

### Difference in $L_e$ of coding and non-coding parts mostly, but not all, caused by difference in p





## Universality class strengthens with increasing genome length (or statistics)

#### Fraction of k-mers whose P-value is less than P3, P6, or P8

	$k = 2 (L_e = 310 b)$			$k=9 \ (L_e=194 \ \text{kb})$		
Length (Mb)	P < P <sub>3</sub>	P <p<sub>6</p<sub>	P <p<sub>8</p<sub>	P <p<sub>3</p<sub>	P <p<sub>6</p<sub>	P <p<sub>8</p<sub>
0.8	0.953	0.906	0.875	0.139	0.0031	0.0001
4.6	0.980	0.960	0.955	0.538	0.418	0.100
30	0.992	0.985	0.979	0.809	0.628	0.519
226	0.997	0.994	0.992	0.930	0.860	0.815

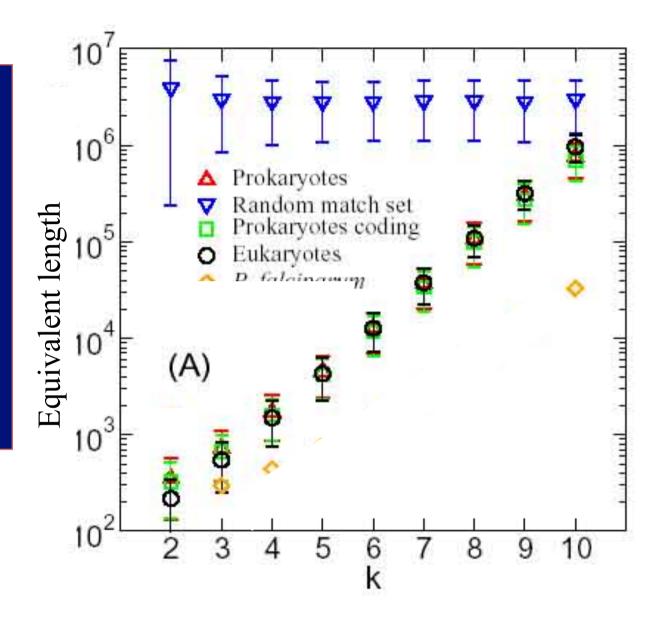
P-values for k-mer distribution given by Eq. (1) (at p = 0.5). Null theory assumes genomes are random sequences. The P-values  $P_3 = 2.7 \times 10^{-3}$ ,  $P_6 = 2.0 \times 10^{-9}$ , and  $P_8 = 1.3 \times 10^{-15}$  correspond to z-values of three, six and eight, respectively. doi:10.1371/journal.pone.0009844.t004

Universality Class of genomic equivalent length

 $\log L_e(k) \sim a k + B$ 

a = 0.40B = 1.69 +- 0.28

Mild exception: Plasmodium



# Is universality class a consequence of similarity in genomes?

Similarity index and similarity matrix

Given a pair of equal-length sequences  $\alpha$  and  $\beta$ , the similarity index  $\eta_{sim}(\alpha,\beta)$  for the pair is defined as

$$\eta_{sim}^{2}(\alpha,\beta) = \frac{1}{k+1} \sum_{m} \frac{1}{2\tau_{m}} \sum_{u \in S_{m}} \frac{(f_{u}^{\{\alpha\}} - f_{u}^{\{\beta\}})^{2}}{\sigma_{m}^{\{\alpha\}} \sigma_{m}^{\{\beta\}}}$$
(15)

where  $S_m$  is an m-set and  $\sigma_m^2$  is the variance of the frequency of the k-mers in  $S_m$ . The pair are similar (in k-mer-content) when  $\eta_{sim} \ll 1$ , are (considered to be) identical when  $\eta_{sim} = 0$ , and are highly dissimilar when  $\eta_{sim} \gtrsim 1$ .

#### Homogeneity of Chromosome on a scale of 10 kb

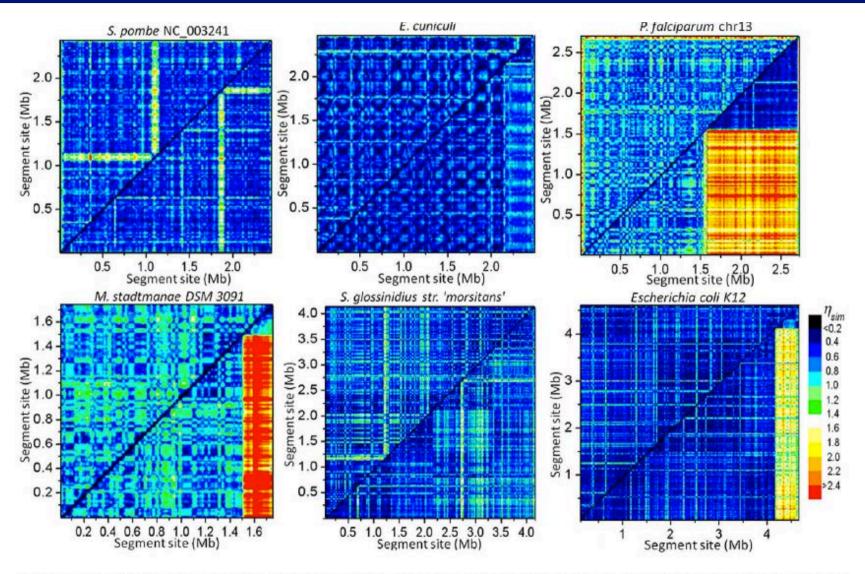
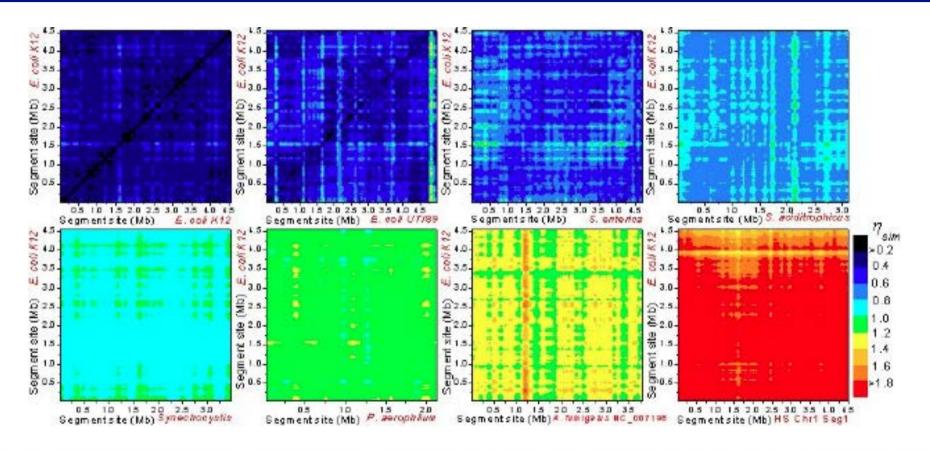


Figure 5. Intra-chromosomes similarity plots. Plots are for k = 2 (Methods). Sliding window has width 25 kb and slide 10 kb; pixel size is 10 kb by 10 kb. In each plot, the coordinates for the upper-left triangle are sites along the chromosome (*chr*), and those for the lower-right triangle are along a concatenate composed of gene (gn, left side) and intergene (ig, right side) parts. In effect, the upper-left triangle shows *chr-chr* similarity, and the lower-right triangle shows gn-gn (lower-left sub-triangle), ig-ig (upper-right sub-triangle), and gn-ig (rectangular) similarities in three separate regions. The lengths of the gn and ig parts are given in Table 3. doi:10.1371/journal.pone.0009844.g005

# Universality class not due to similarity in genomes



**Figure 6. Intra-** *E. coli* **and inter-chromosome similarity plots.** The plots are those of *E. coli* chromosome *vs.* the chromosomes of, left to right and top to bottom, *E. coli*, *E. coli UT189, Salmonella*, the delta-proteobacteria *S. aciditrophicus*, the cyanobacteria *Synechocystis*, the archaea *P. aerophilum*, chromosome 5 of the fungus *A. fumigatus*, and the first 4.5 Mb segment from chromosome 1 of *H. sapiens*. Coordinates are sites along the sequence. Sliding window width is 100 kb and slide is 25 kb, pixel size is 25 kb by 25 kb. doi:10.1371/journal.pone.0009844.g006

### Summary on genomic Le

- Universality class for fixed word length k, L<sub>e</sub> is (approximately) the same for all genomes
  - Log  $L_e(k) = ak + B$ ; a, B universal constants (for ~900 complete sequence)
- For k < 8,  $L_e$  is much shorter than true genome length;  $L_e$  (k=2) ~ 300 b
- L<sub>e</sub> for coding and non-coding parts about the same
- Universality not due sequence similarity

# Small $L_e$ is a signature of segmental duplication

- Recall:  $L_e$  is the length of random sequence having genomic CV
- Take a random sequence of length  $L_{\rm e}$  and replicate it n times, then sequence length is  $nL_{\rm e}$  but equivalent length is still  $L_{\rm e}$  (for all  $k << L_{\rm e}$ )
- Hint: small genomic  $L_e$  caused by segmental duplication

#### Fun with $l_e$ and concatenates

$$I_e = \tau/(CV)^2 = \tau/(\sigma/f)^2 = L^2 \tau^3/\sigma^2 \quad (\tau = 4k)$$

Let X be genomic, and R random, sequence of same length;

RX be concatenate of R and X.

Then 
$$\sigma(X) >> \sigma(R)$$
,  $I_e(X) << I_e(R) = L(R) = L(X)$ .

$$I_{e}(RX) = L(RX)^{2} \tau^{3}/\sigma(RX)^{2} \sim (L(R) + L(X))^{2} \tau^{3}/\sigma(X)^{2}$$
$$= 4L(X))^{2} \tau^{3}/\sigma(X)^{2} = 4 I_{e}(X)$$

Similarly,

$$I_e(RXX) \sim (3/2)^2 I_e(X);$$
  $I_e(RRX) \sim 9 I_e(X)$ 

Similar (notation ~): similar in word content. Suppose X and X' are same length sequences from the same genome, Y from a different genome, and R is a random sequence. XY denote concatenate of X and Y.

- $-I_e(X) \sim I_e(X') \sim I_e(XX')$
- If X~Y, then  $I_e(X) \sim I_e(Y)$ ,  $I_e(XY) \sim I_e(X)$ ,
- If X not~ Y, then  $I_e(XY) > \min(I_e(x), I_e(Y))$
- $I_e(RX)$  approx 4  $I_e(X)$
- $I_e(RXY)$  approx 2.3  $I_e(XY)$
- $-I_e(RR'X)$  approx  $9I_e(X)$

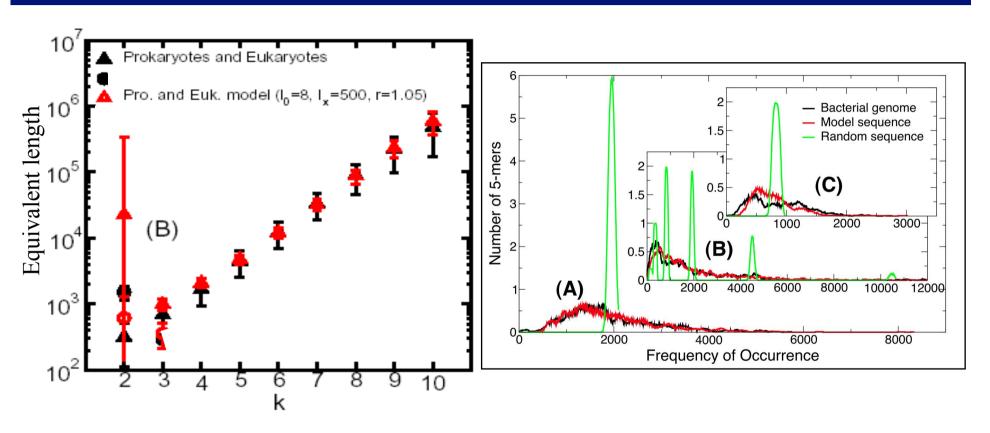
	L <sub>e</sub>						
	k = 2		k=6				
Sequence	<i>l</i> = 50	<i>l</i> = 200	<i>l</i> = 50	<i>l</i> = 200			
R	47.5 ± 28.2	154±126	48.6 ± 1.5	192±5			
RR'	$37.0 \pm 16.2$	124 ± 46	$48.2\pm1.2$	197±5			
Α	$.348 \pm .037$	$.360 \pm .033$	$9.55\pm.69$	11.7 ± .7			
AA'	$.357 \pm .046$	$.352 \pm .023$	$9.88 \pm 1.07$	11.1 ± .7			
AC <sub>1</sub>	$.351 \pm .061$	$.361 \pm .021$	$9.37 \pm 1.01$	11.5 ± .6			
AC <sub>2</sub>	$.354 \pm .043$	$\textbf{.384} \pm \textbf{.045}$	$9.18\pm.83$	11.6 ± .9			
AC <sub>3</sub>	$.359 \pm .051$	$.371 \pm .034$	11.0 ± .9	$14.2 \pm 1.5$			
AD <sub>1</sub>	.411±.044	$\textbf{.423} \pm \textbf{.024}$	$11.8 \pm .9$	$14.3 \pm .6$			
AD <sub>2</sub>	.942±.275	$1.05\pm.09$	14.9 ± 1.4	$20.4 \pm 1.1$			
AD <sub>3</sub>	$.598 \pm .104$	$.613 \pm .052$	$17.9 \pm 1.6$	$24.0 \pm 1.6$			
AD <sub>4</sub>	$.324 \pm .052$	$\textbf{.383} \pm \textbf{.055}$	11.2 ± 1.9	$16.9 \pm 1.9$			
В	$.124 \pm .029$	$\textbf{.166} \pm \textbf{.099}$	$5.17\pm.68$	$6.54 \pm 2.00$			
BB'	.232±.155	$.258 \pm .183$	$6.16 \pm 1.94$	$7.54 \pm 2.30$			
AB	.463 ± .241	.502 ± .263	11.2 ± 1.9	15.2 ± 3.5			
RA	1.19±.09	1.34 ± .20	22.6 ± 1.2	38.5 ± 3.0			
RB	.575±.321	.754±.637	15.6 ± 4.2	23.3 ± 8.5			
RAB	873±.424	1.10±.49	18.4 ± 3.2	31.3 ± 6.0			
RR' A	2.63±.66	3.16±.30	31.5 ± 2.1	72.2 ± 6.8			
RR'B	1.03±.62	1.37 ± .70	22.9 ± 4.5	44.7 ± 14.3			

### $l_e$ of a sequence after n-fold growth

Equivalent length  $(l_e)$  of a sequence after an n-fold increase in length via three basic modes of growth. Initial sequence length is  $l_0 >> 1$ , final sequence length is  $L = nl_0$ .

Sequence type	Initial $l_e$	Mode of growth	Final $l_e$
Random	$l_0$	Random base-by-base growth	$L(=nl_0)$
		Whole-sequence replication $(n-1 \text{ times})$	$pprox l_0$
		Segmental duplication	$l_0 < l_e << nl_0$
Non-random	$l_{e0} \ (<< l_0)$	Random base-by-base growth	$\approx \min(n^2 l_{e0}, L)$
		Whole-sequence replication $(n-1 \text{ times})$	$pprox l_{e0}$
		Segmental duplication	$l_{e0} < l_e < < \min(n^2 l_{e0}, L)$

## RSD model with three universal parameters generates artificial genomes with universal $L_e$



Note: Deep consequences in understanding speed of evolution (Lecture 4)

HD Chen, et al. Universal Global Imprints of Genome Growth and Evolution - Equivalent Length and Cumulative Mutation Density. PLoS ONE (2010) 5(4): e9844.

### End of Lecture One