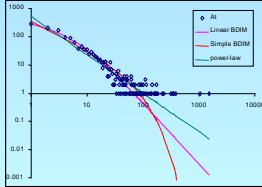



National Center for Biotechnology Information

## Birth, Death and Innovation Models (BDIM) of genome evolution

Eugene V. Koonin  
NCBI, NLM, NIH, Bethesda, MD





$$df_i(t)/dt = +\lambda_{i-1}f_{i-1} - \delta f_i - \lambda f_i + \delta_{i+1}f_{i+1}$$

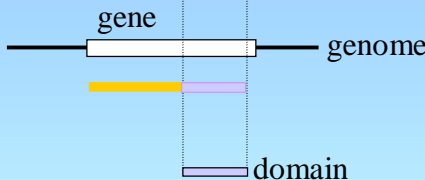
National Center for Biotechnology Information

- mathematical formulation and classification of models
- where do the power laws come from
- empirical data fit
- do BDIM tell us anything biologically non-trivial?

National Center for Biotechnology Information

## BDIM: basic concepts

part-of-gene-encoding-  
an-individual-domain = "domain"

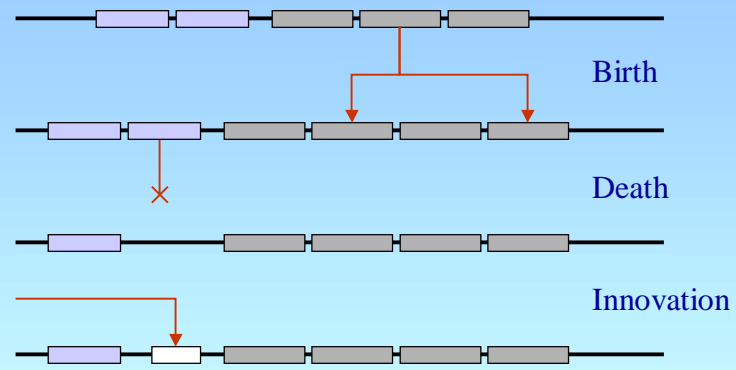


The diagram illustrates the relationship between a gene, a genome, and a domain. A horizontal line represents the genome, with a white box labeled 'gene' on it. Below the gene box, a yellow and pink bar represents the gene's structure. A pink box labeled 'domain' is shown below the pink part of the gene bar, with vertical dashed lines indicating its position within the gene.

Genome: a finite "bag" of independently evolving domains.  
Domains form families of paralogs (domain families for short).

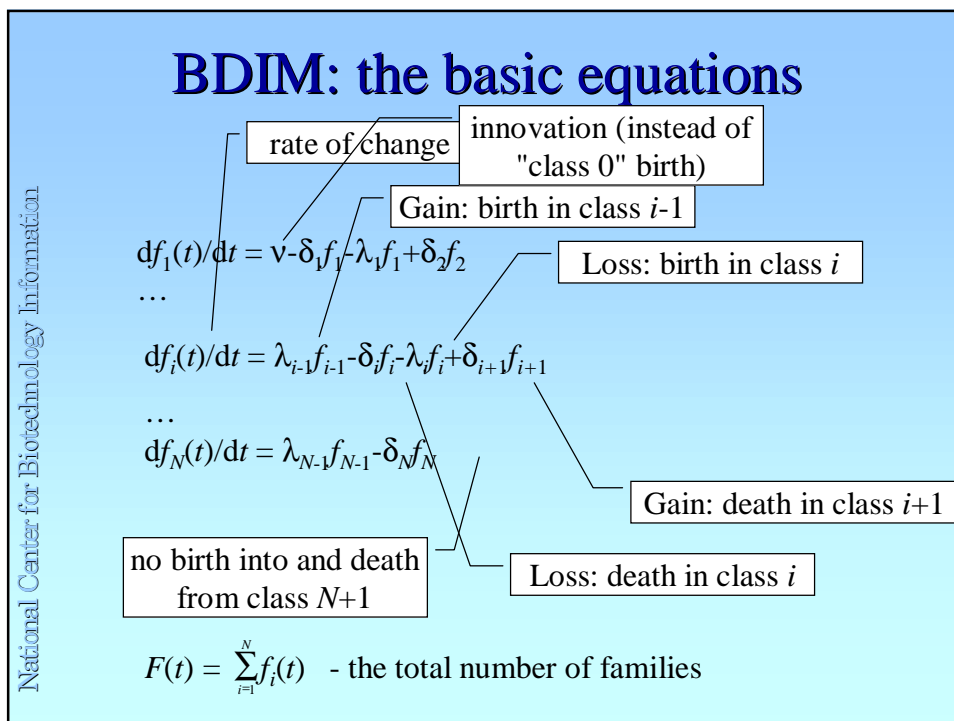
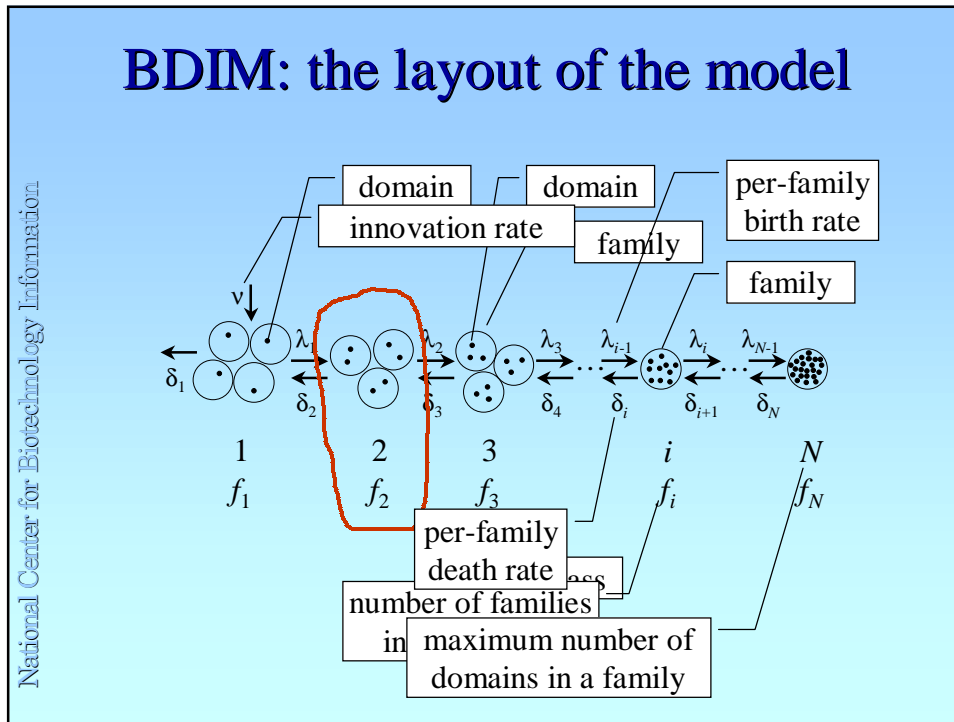
National Center for Biotechnology Information

## BDIM: elementary events



The diagram shows three horizontal lines representing genes. The top line has five domains (two purple, three grey). Red arrows point to the second and third domains, labeled 'Birth'. The middle line has the same five domains, but the second domain is crossed out with a red 'X', labeled 'Death'. The bottom line has the same five domains, but the second domain is white, labeled 'Innovation'.

BDIM – Birth, Death and Innovation Model



National Center for Biotechnology Information

## BDIM: the equilibrium

$df_i(t)/dt = 0$  - equilibrium for the number of domain families in each size class  
 $dF(t)/dt = 0$  - equilibrium for the total number of families

- There exists a **unique** and **stable** equilibrium state  $f_1, f_2, \dots, f_N$
- The model reaches equilibrium **exponentially**:  
 $|f_i(t) - f_i| \sim e^{-kt}$
- The model is "open" at one end only (class 1 families).  
 A simple condition describes the equilibrium for the **total** number of families:

$v = \delta_1 f_1$   
 innovation Death of families

National Center for Biotechnology Information

## Simple BDIM

**Independence hypothesis:**

- all elementary events are independent of each other
- the rates of **individual domain** birth ( $\lambda$ ) and death ( $\delta$ ) do not depend on  $i$  (number of domains in a family).

**Corollary:**

$$\lambda_i = \lambda i$$

$$\delta_i = \delta i$$

**The basic equation for domain family evolution:**

$$df_i(t)/dt = \lambda(i-1)f_{i-1} - (\delta + \lambda)if_i + \delta(i+1)f_{i+1}$$

National Center for Biotechnology Information

## Simple BDIM: equilibrium

The equilibrium solution:  
 $f_i \sim (\lambda/\delta)^i/i$  - truncated logarithmic distribution  
 If  $\lambda = \delta$   
 $f_i \sim 1/i$  - power law (degree = -1)

National Center for Biotechnology Information

## BDIM hierarchy

Master BDIM

 $\lambda_i = \lambda(i)$   
 $\delta_i = \delta(i)$

↓

rational

 $\lambda_i = P(i)/Q(i)$

↓

polynomial

 $\lambda_i = p_0 + p_1i + p_2i^2 + p_3i^3 \dots$

↓

simple

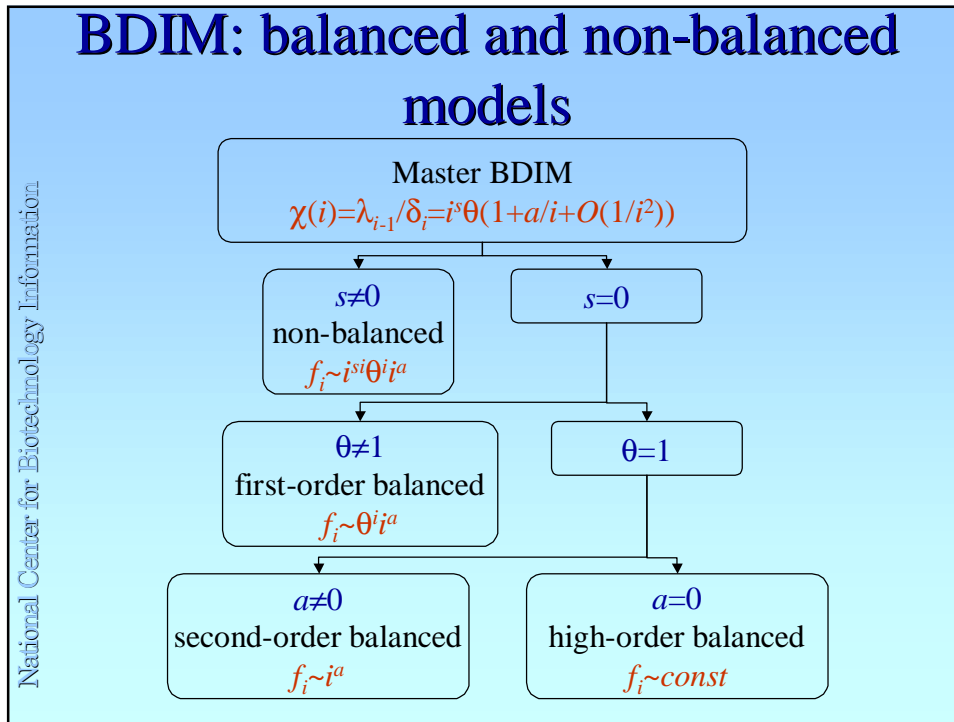
 $\lambda_i = \lambda i$

linear

 $\lambda_i = \lambda(i+a)$

quadratic

 $\lambda_i = \lambda(i+a)(i+a_1)$

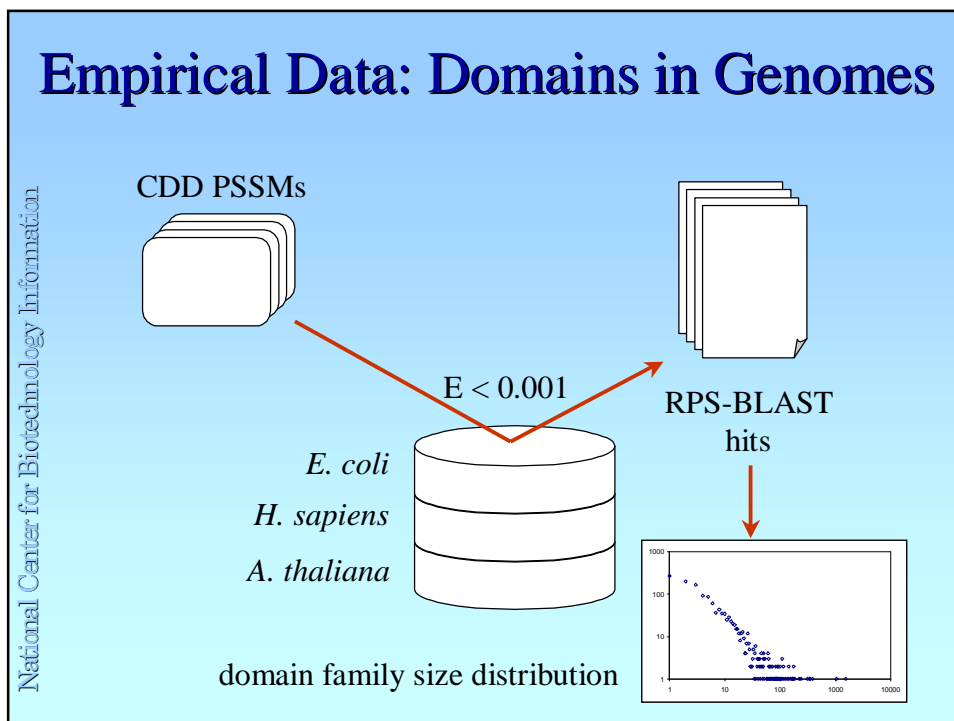
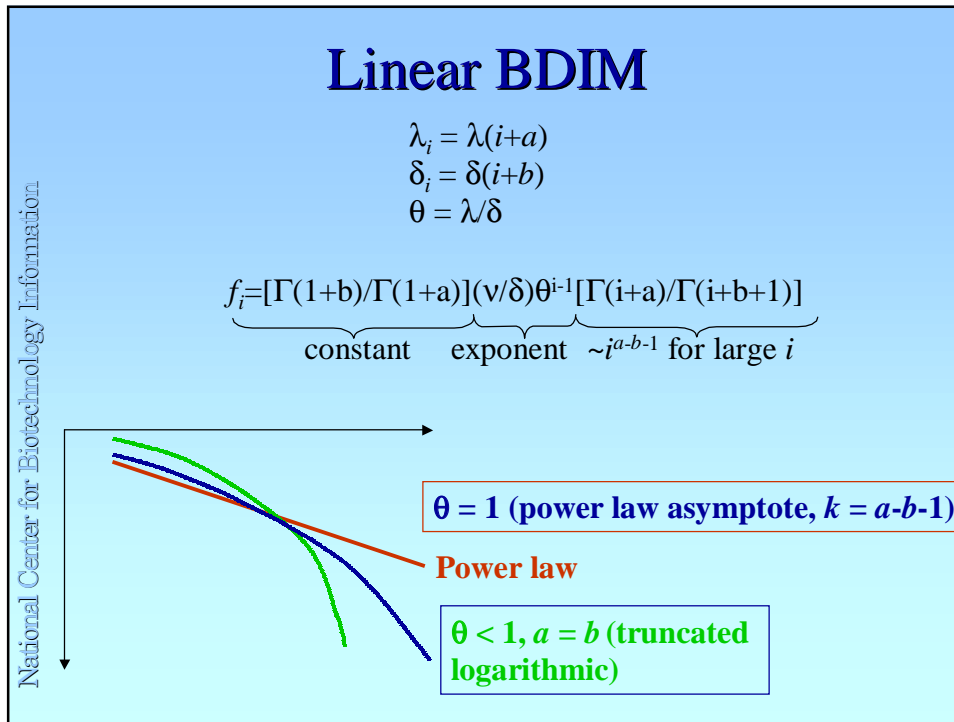


## BDIM: only first/second order balanced models make sense

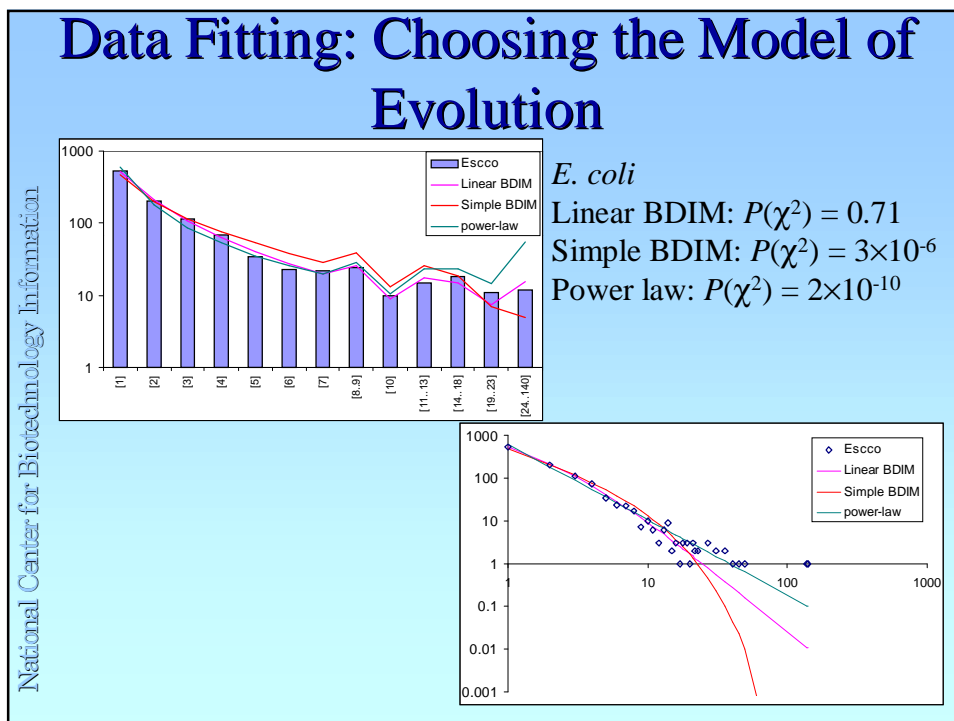
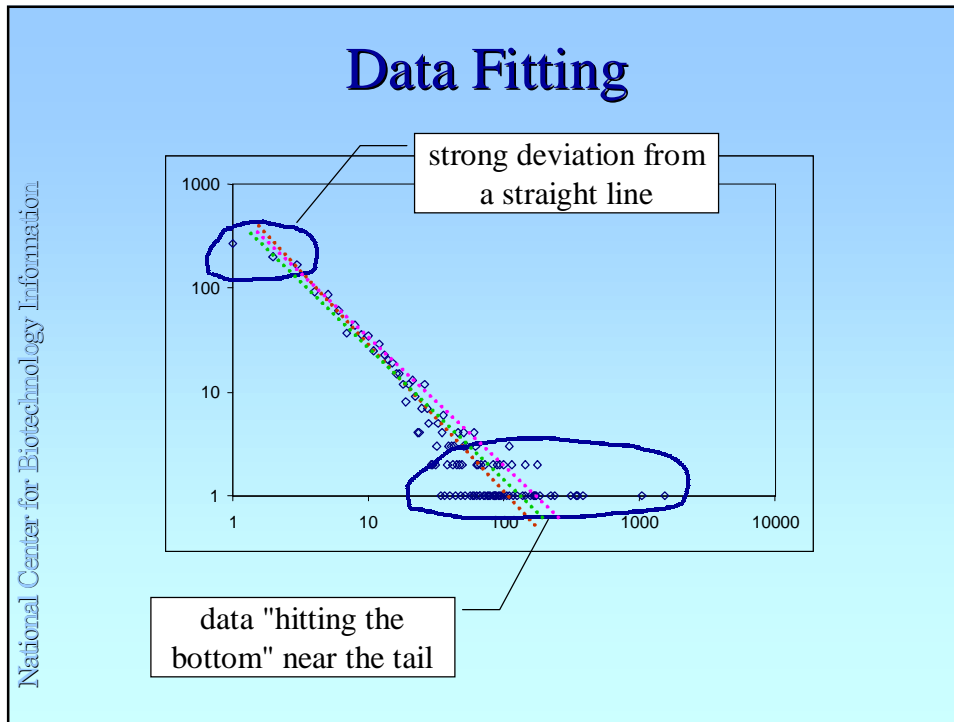
National Center for Biotechnology Information

- Non-balanced BDIM: unrealistic family size distributions with extremely strong dependence on  $i$ :  
**either no large families at all or mostly large families**
- High-order balanced BDIM: equally unrealistic, uniform distribution

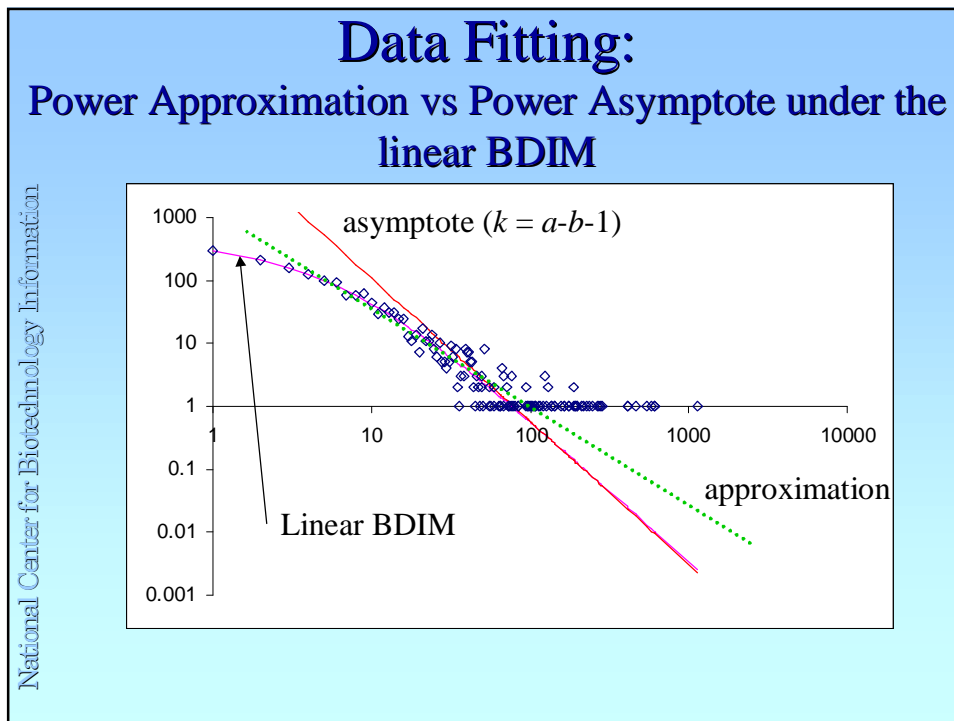
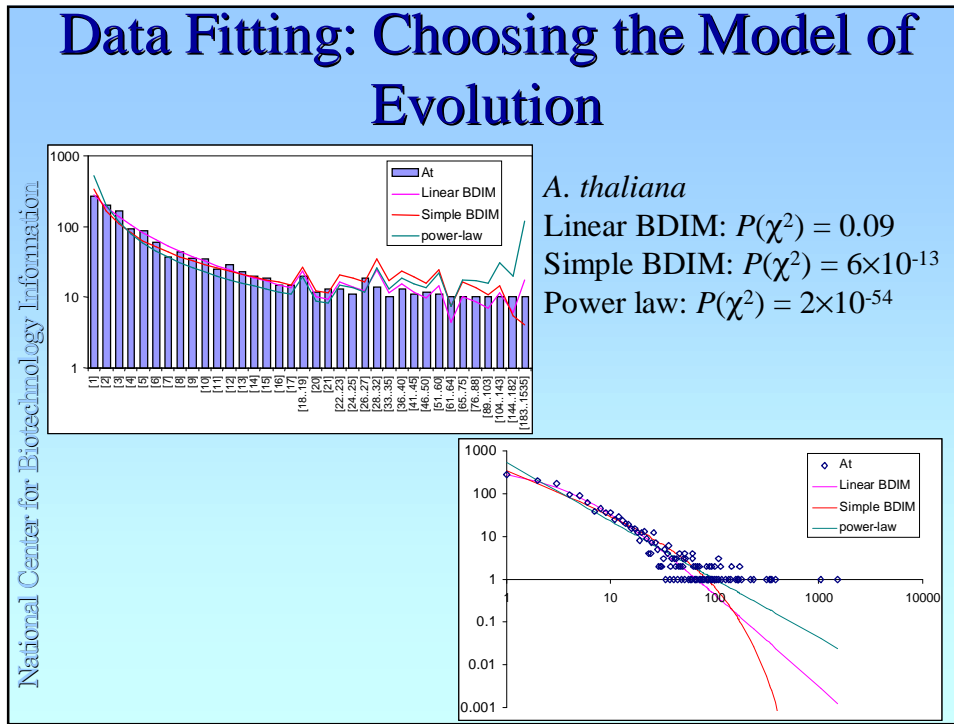
Birth and Death of Domains: Simple Models of Genome Evolution Explain Power Law Distribution of Protein Family Size



# Birth and Death of Domains: Simple Models of Genome Evolution Explain Power Law Distribution of Protein Family Size



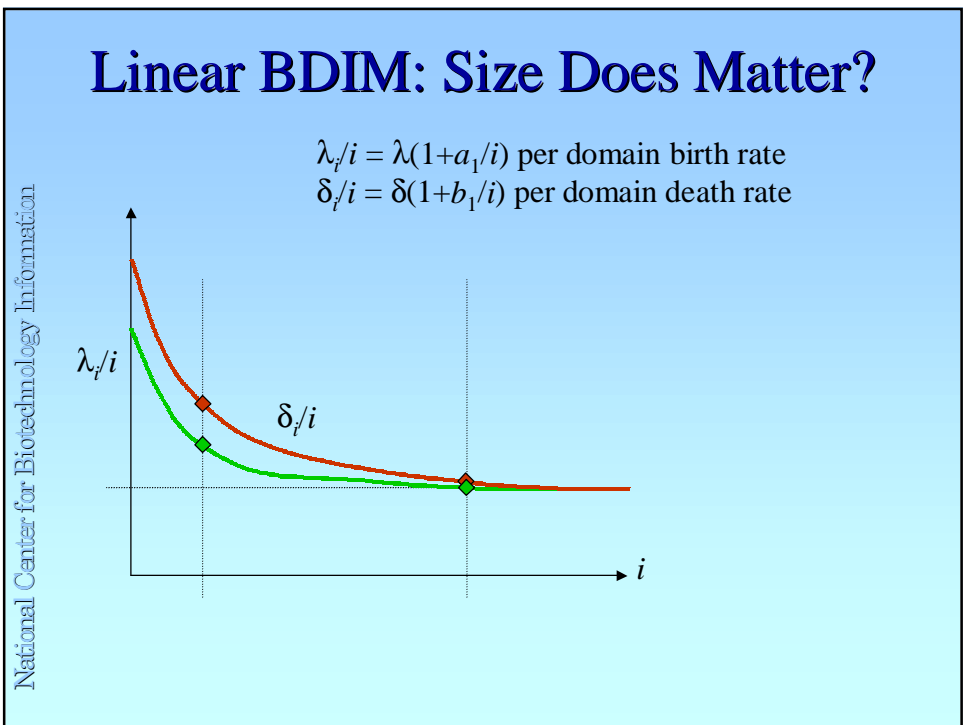




## Linear BDIM parameters and other relevant numbers for bacteria, archaea and eukaryotes

	nORF	nFam	nDom	nHit	iMax	$d_1 (d'_1)$	$a$	$b$	$k$	$v/\delta=v/\lambda$	$G$
Sc	6340	1080	4575	3331	130	420 (436)	1.55	3.27	-2.72	1861.8	[3.28..3.53]
Dm	13605	1405	11734	7262	335	426 (435)	1.62	2.79	-2.17	1648.2	[8.44..15.50]
Ce	20524	1418	17054	11090	662	423 (421)	1.13	2.03	-1.89	1273.0	[16.18..∞]
At	25854	1405	21238	15006	1535	270 (277)	3.80	4.98	-2.18	1657.7	[17.09..26.89]
Hs	39883	1681	27844	16755	1151	298 (288)	5.16	6.43	-2.27	2136.2	[17.14..22.88]
Thema	1846	772	1683	1268	97	501 (499)	0.14	2.22	-3.08	1606.4	[1.04..1.06]
Metth	1869	693	1480	1150	43	438 (436)	0.12	2.00	-2.88	1305.3	[1.18..1.27]
Sulso	2977	695	1950	1614	81	386 (385)	0.36	2.04	-2.68	1167.8	[1.83..2.00]
Bacsu	4100	1002	3413	2502	124	507 (510)	0.48	2.01	-2.53	1534.6	[2.46..2.79]
Escco	4289	1078	3624	2765	140	523 (519)	0.84	2.54	-2.70	1837.0	[2.45..2.61]

$G$ =total duplication rate/innovation rate



## Conclusions

- I. Only balanced BDIM produce reasonable equilibrium distributions of domain family size; equilibrium is reached rapidly, suggesting a "punctuated equilibrium"-like mode of genome evolution.**
- II. The simplest evolutionary model that adequately describes the observed distribution of domain family size is the linear, second-order balanced BDIM; accordingly, per-domain birth/death rate depends on family size, the larger families being less dynamic in evolution.**
- III. The rates of domain innovation and birth are comparable.**

National Center for Biotechnology Information

**The original version of BDIM is fully deterministic.**

**In order to be able to explore the dynamics of genome evolution, we introduce a stochastic (Markov) version.**

National Center for Biotechnology Information

National Center for Biotechnology Information

**Markov version of BDIM (0 class introduced)**  
(innovation interpreted as extraction from class 0)

$$\begin{aligned} \frac{d p_0(t)}{dt} &= -\lambda_0 p_0(t) + \delta_1 p_1(t), \\ \frac{d p_1(t)}{dt} &= \lambda_0 p_0(t) - (\lambda_1 + \delta_1) p_1(t) + \delta_2 p_2(t), \\ \dots \\ \frac{d p_i(t)}{dt} &= \lambda_{i-1} p_{i-1}(t) - (\lambda_i + \delta_i) p_i(t) + \delta_{i+1} p_{i+1}(t) \text{ for } 1 < i < N, \\ \dots \\ \frac{d p_N(t)}{dt} &= \lambda_{N-1} p_{N-1}(t) - \delta_N p_N(t) \end{aligned}$$

**Modified Markov version of BDIM (no 0 class, class 1 immortal)**

$$\begin{aligned} \frac{d p_1(t)}{dt} &= -\lambda_1 p_1(t) + \delta_2 p_2(t), \\ \dots \\ \frac{d p_i(t)}{dt} &= \lambda_{i-1} p_{i-1}(t) - (\lambda_i + \delta_i) p_i(t) + \delta_{i+1} p_{i+1}(t) \text{ for } 1 < i < N, \\ \dots \\ \frac{d p_N(t)}{dt} &= \lambda_{N-1} p_{N-1}(t) - \delta_N p_N(t). \end{aligned}$$

no innovation

National Center for Biotechnology Information

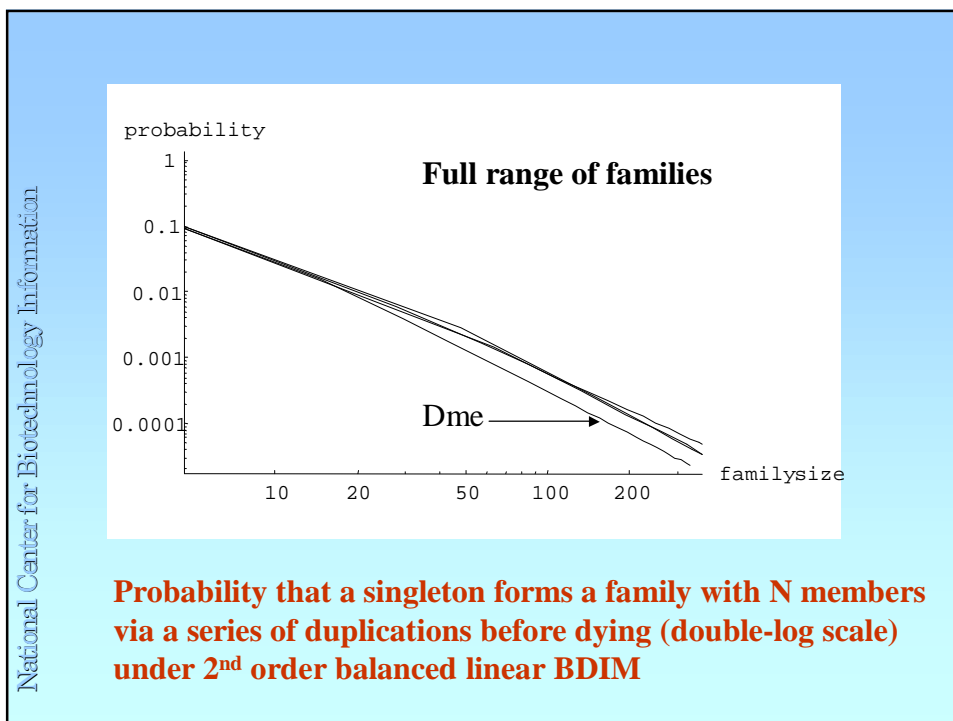
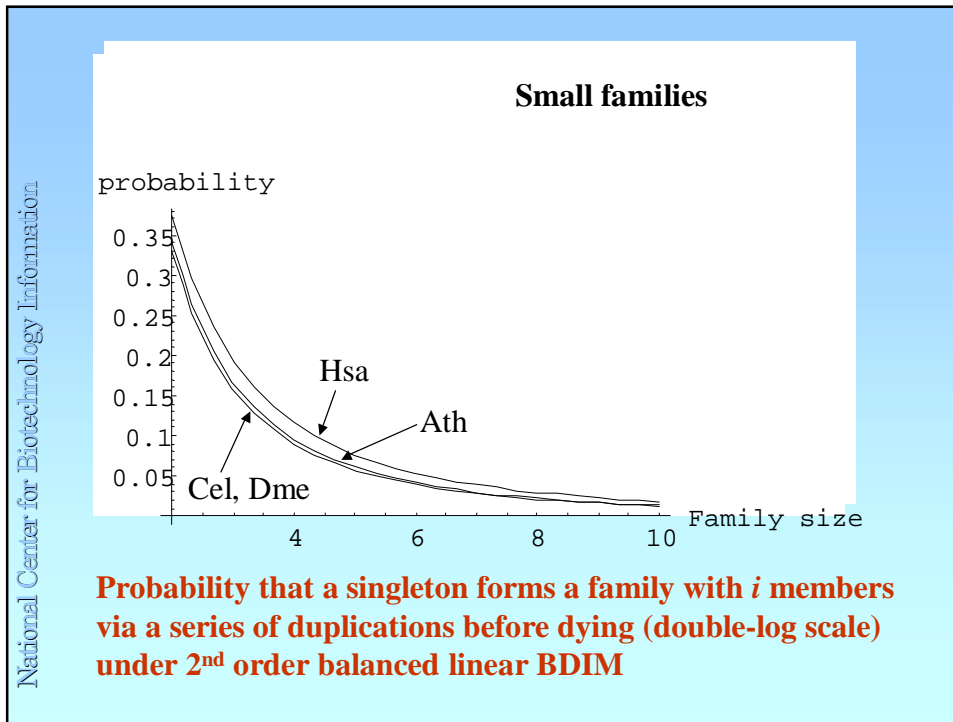
**Probability for a family to reach size  $n$  from size  $i$  before extinction (size 0)**

$$P(i;n) = \left(1 + \sum_{j=1}^{i-1} \prod_{k=1}^j \delta_k / \lambda_k\right) / \left(1 + \sum_{j=1}^{n-1} \prod_{k=1}^j \delta_k / \lambda_k\right)$$

And, for 2<sup>nd</sup> order balanced linear BDIM and  $i=1$ ,

$$P(1,n) = 1 / \left(1 + \frac{\Gamma(1+a)}{\gamma \Gamma(1+b)} \left( \frac{\Gamma(b+n+1)}{\Gamma(a+n)} - \frac{\Gamma(2+b)}{\Gamma(1+a)} \right)\right)$$

Birth and Death of Domains: Simple Models of Genome Evolution Explain Power Law Distribution of Protein Family Size

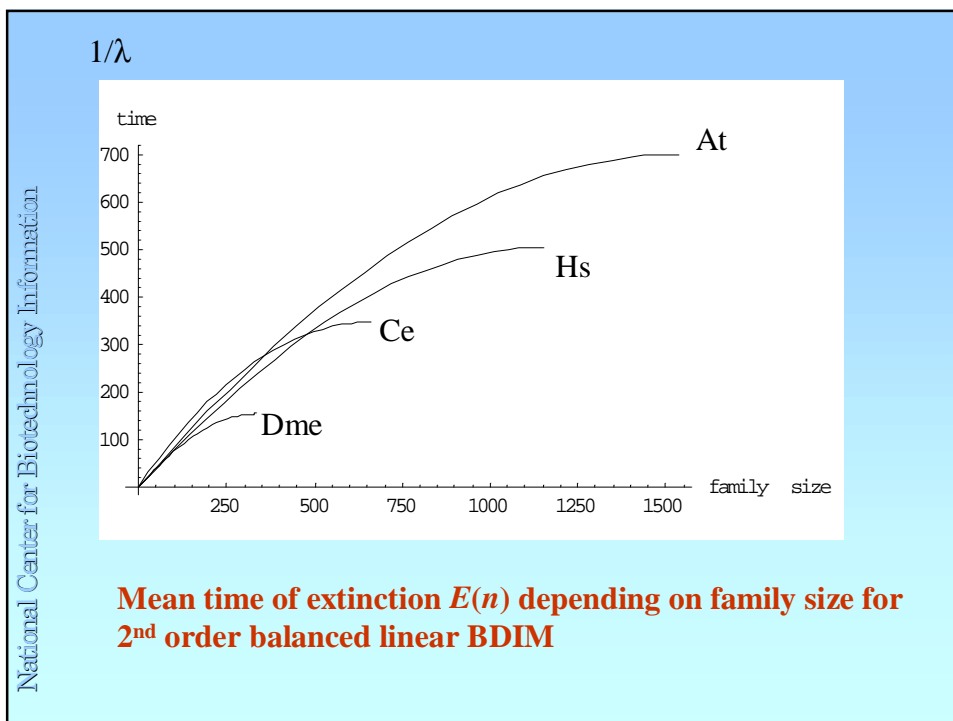


National Center for Biotechnology Information

**Time before extinction of a family of size  $n$  for 2<sup>nd</sup> order balanced linear BDIM**

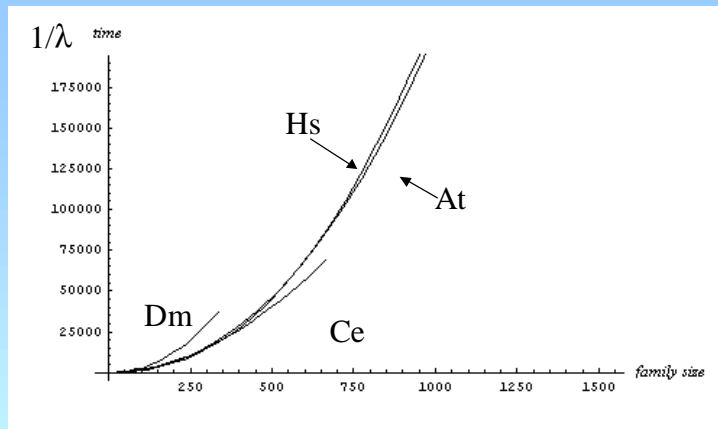
$$E(n) = 1/\lambda \sum_{j=1}^n \frac{\Gamma(j+b)}{\Gamma(j+a)} \sum_{k=j}^N \frac{\Gamma(k+a)}{\Gamma(k+1+b)}$$

**For the time being, we use  $1/\lambda$  as a natural time scale for BDIM...**



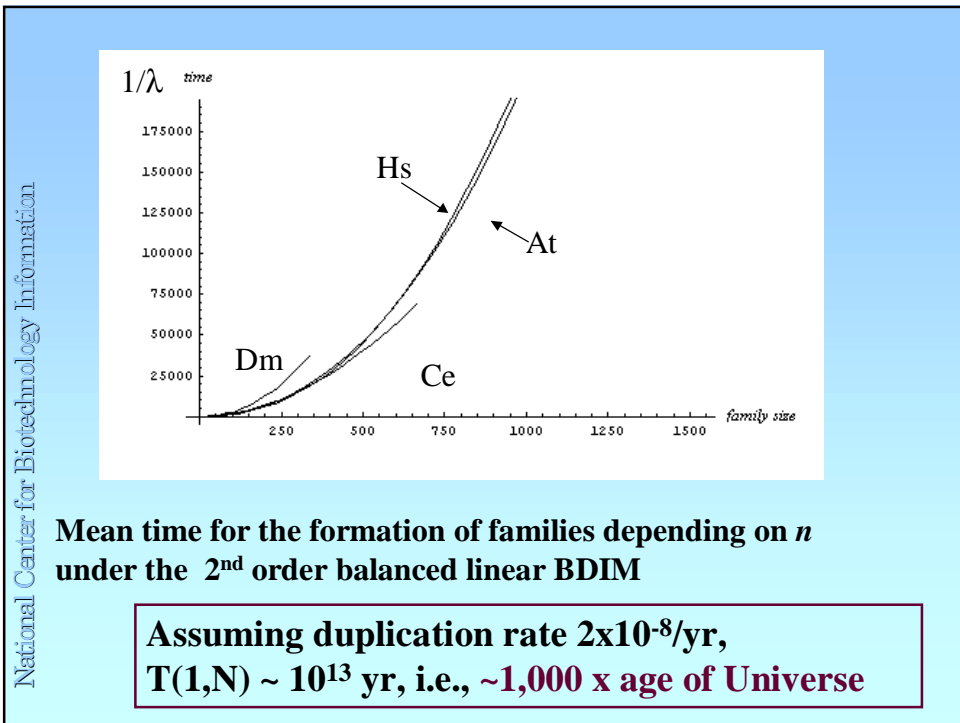
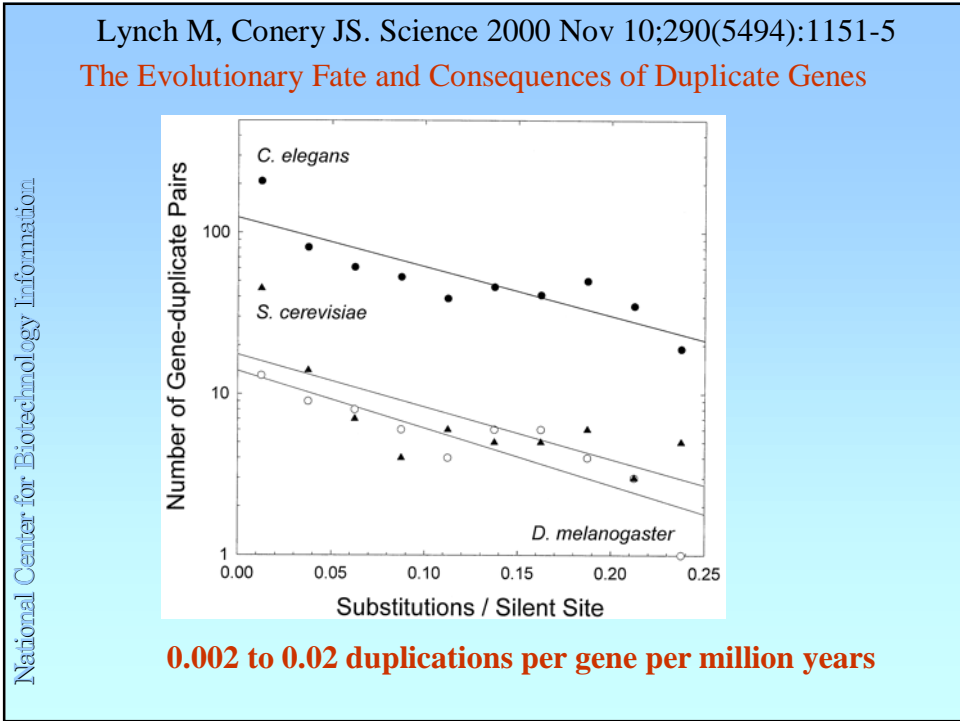
**Time required for a singleton to reach family size  $n$  under 2<sup>nd</sup> order balanced linear BDIM (measured in  $1/\lambda$  units)**

$$M(1;n) = 1/\lambda \sum_{k=1}^{n-1} \left( \frac{\Gamma(b+k+1)}{\Gamma(a+k+1)} \sum_{i=1}^k \frac{\Gamma(a+i)}{\Gamma(b+i+1)} \right). \quad (8.7)$$



**Mean time for the formation of families depending on  $n$  under the 2<sup>nd</sup> order balanced linear BDIM**

Birth and Death of Domains: Simple Models of Genome Evolution Explain Power Law Distribution of Protein Family Size





National Center for Biotechnology Information

**Stochastic characteristics of the linear BDIM:**

- i) extremely large difference between times of formation and extinction of the largest families for some genomes – extinction happens **much** faster;
- ii) Under the available empirical estimates of duplication rate the time required for the formation of the largest families is **unrealistically** long

National Center for Biotechnology Information

Consequently, we must **replace the linear BDIM** with another model such that:

- 1) the stationary distribution of the family sizes is the same as for the linear BDIM;
- 2) the new model provides for much more rapid evolution of gene families under realistic values of duplication and deletion rates;
- 3) the ratio of family formation and extinction mean times must be significantly less than for the linear BDIM.

To obtain a new BDIM without changing the stationary distribution:

$$\lambda_i^* \rightarrow \lambda_i g(i), \quad \delta_i^* \rightarrow \delta_i g(i-1)$$

$$g > 0 \\ g(0) = 1$$

Probability of formation of family of size  $n$  prior to extinction:

$$P^*(1, n) = 1 / \left( 1 + \frac{\Gamma(1+a)}{\Gamma(1+b)} \sum_{k=1}^{n-1} \frac{1}{g(k)} \frac{\Gamma(b+k+1)}{\Gamma(a+k+1)} \right);$$

Extinction time:

$$E_s^* = 1 / \lambda \sum_{k=1}^s \sum_{i=k}^N 1/g(k-1) \left[ \frac{\Gamma(a+i)}{\Gamma(b+i+1)} \frac{\Gamma(b+k)}{\Gamma(a+k)} \right]$$

Family formation time (class 1 immortal):

$$M^*(1; n) = 1 / \lambda \sum_{k=1}^{n-1} \frac{1}{g(k)} \frac{\Gamma(b+k+1)}{\Gamma(a+k+1)} \sum_{i=1}^k \frac{\Gamma(a+i)}{\Gamma(b+i+1)}.$$

National Center for Biotechnology Information

### Quadratic BDIM

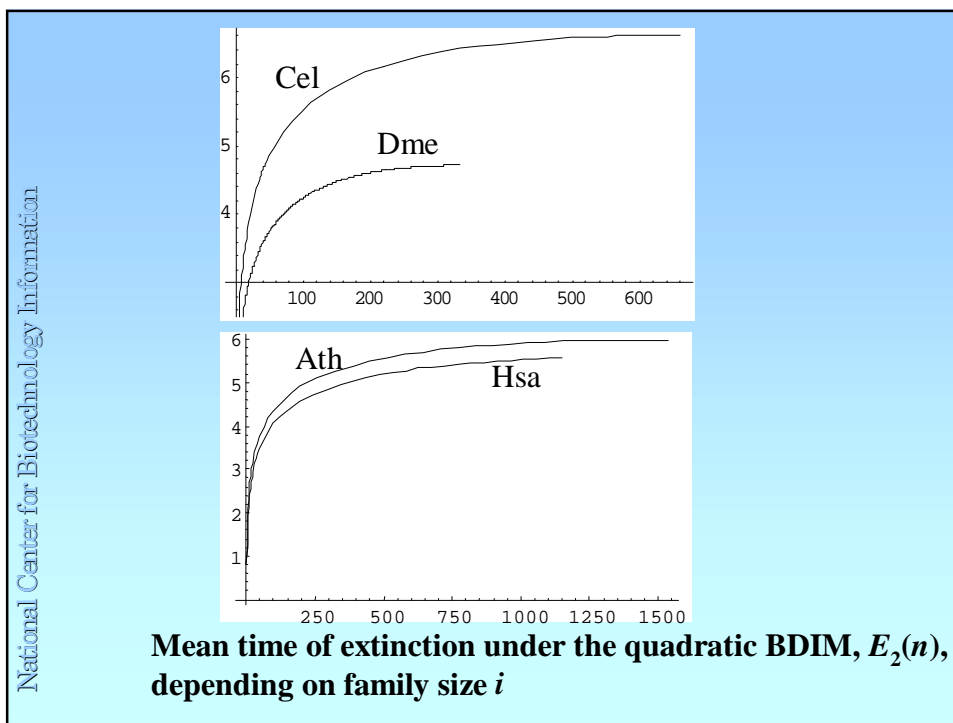
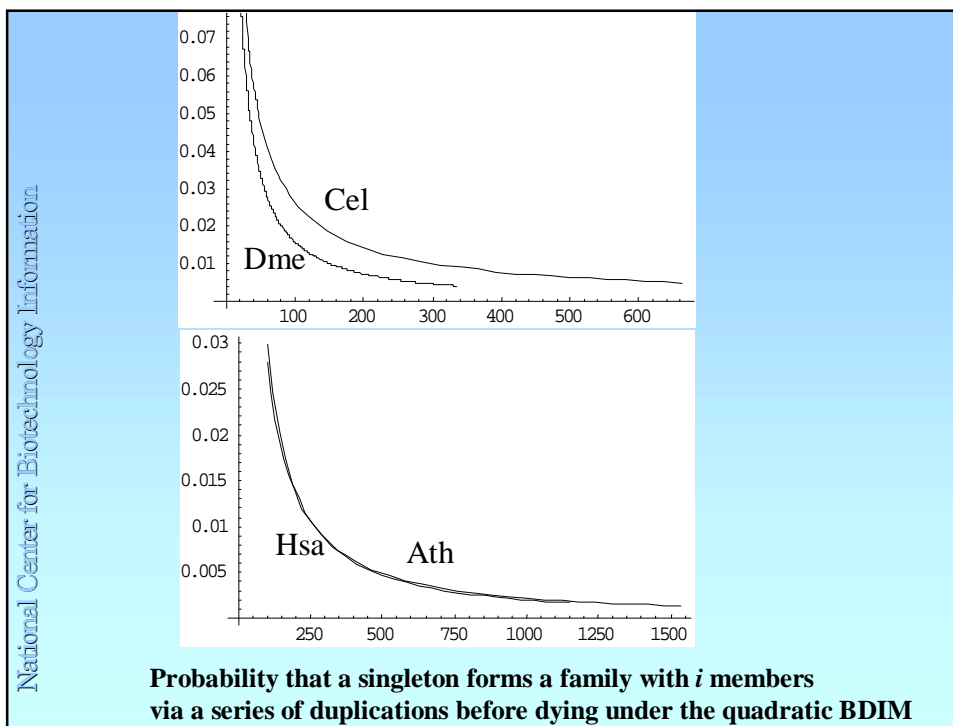
Let us consider the 2<sup>nd</sup> order balanced polynomial (quadratic) Markov BDIM with duplication and deletion rates

$$\lambda_i = \lambda(i+a)(i+1), \\ \delta_i = \lambda(i+b)i.$$

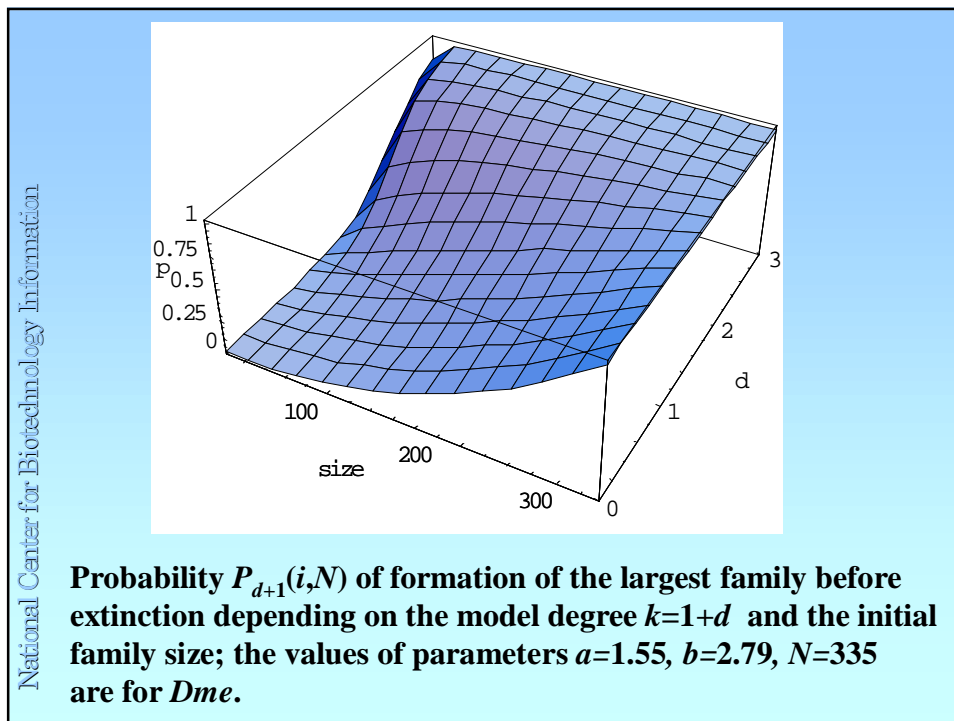
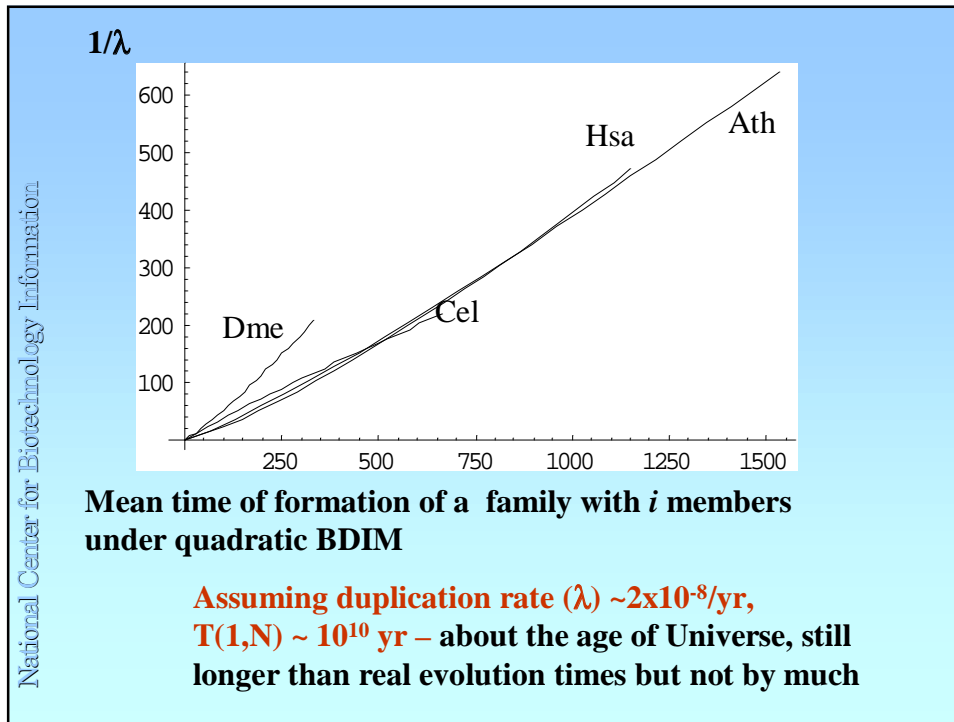
This may be interpreted as introducing pairwise interactions between family members

National Center for Biotechnology Information

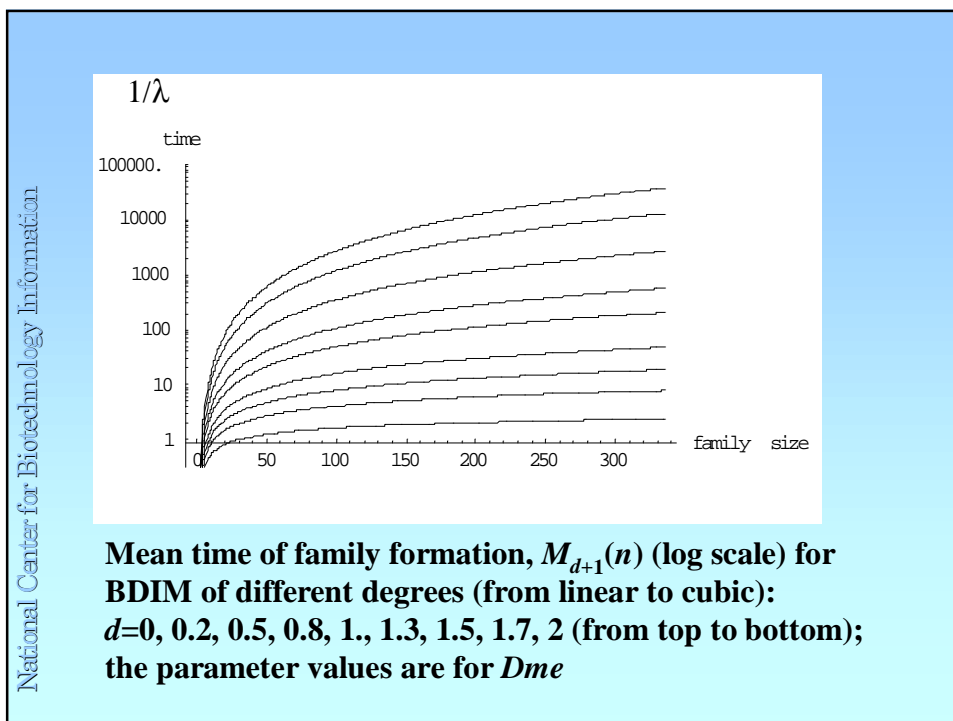
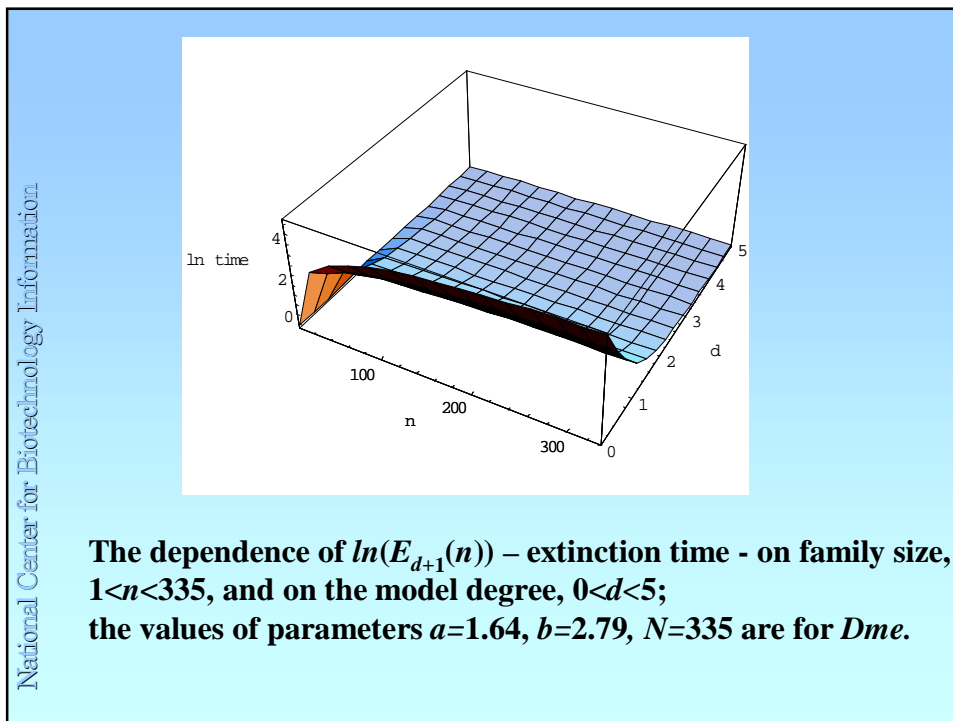
# Birth and Death of Domains: Simple Models of Genome Evolution Explain Power Law Distribution of Protein Family Size



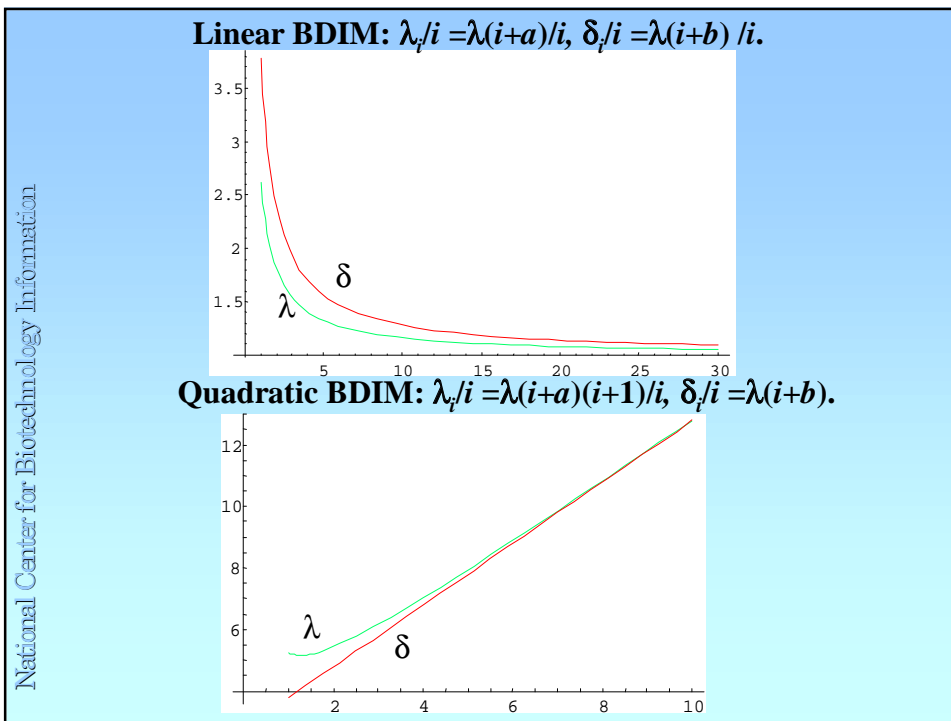
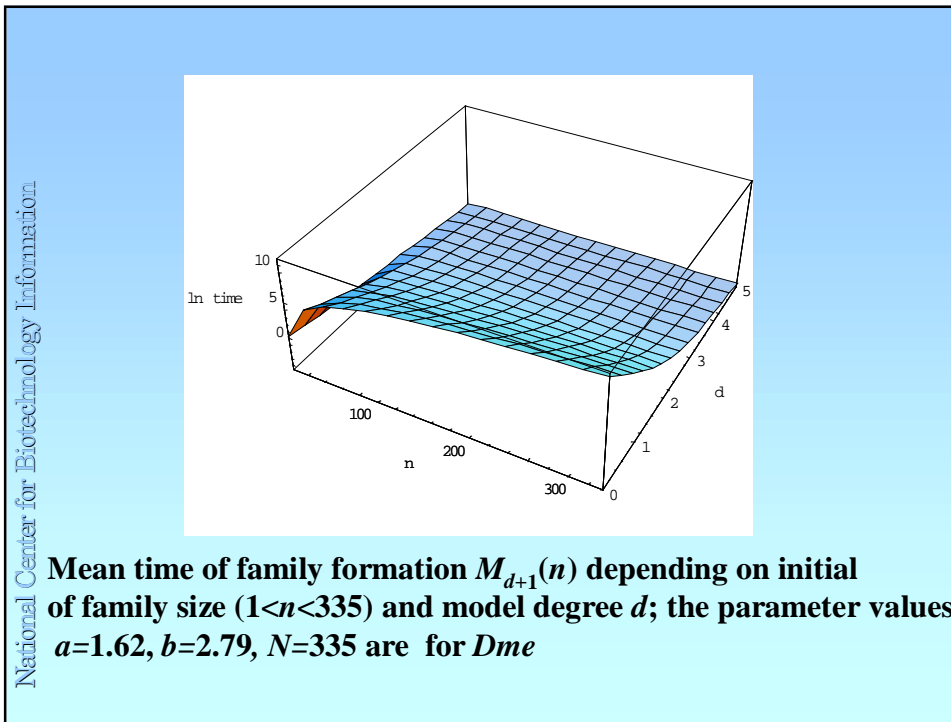
Birth and Death of Domains: Simple Models of Genome Evolution Explain Power Law Distribution of Protein Family Size



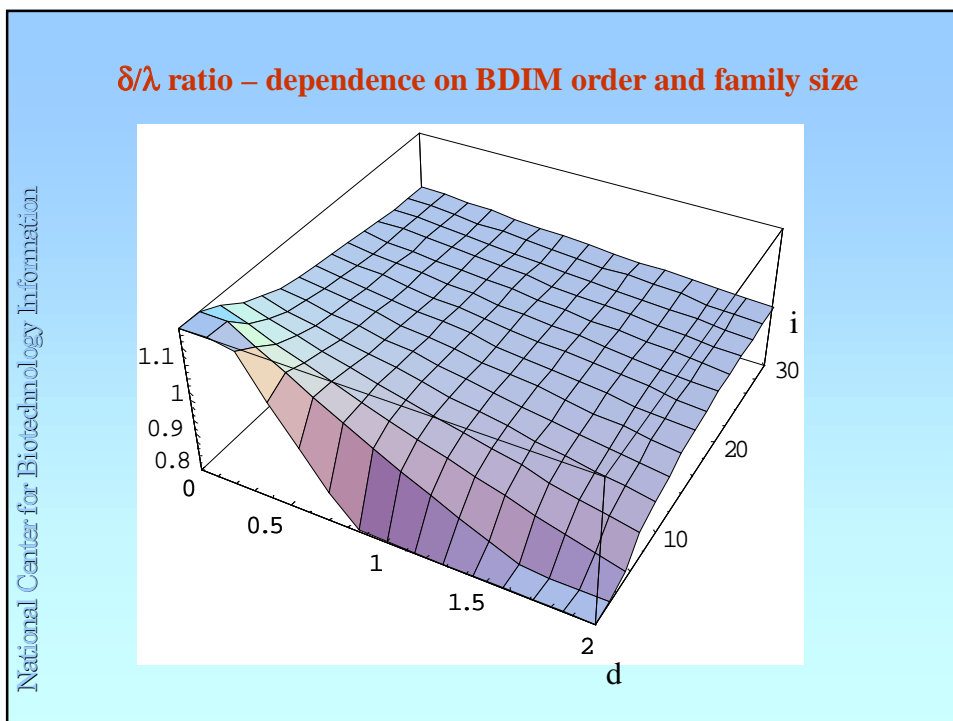
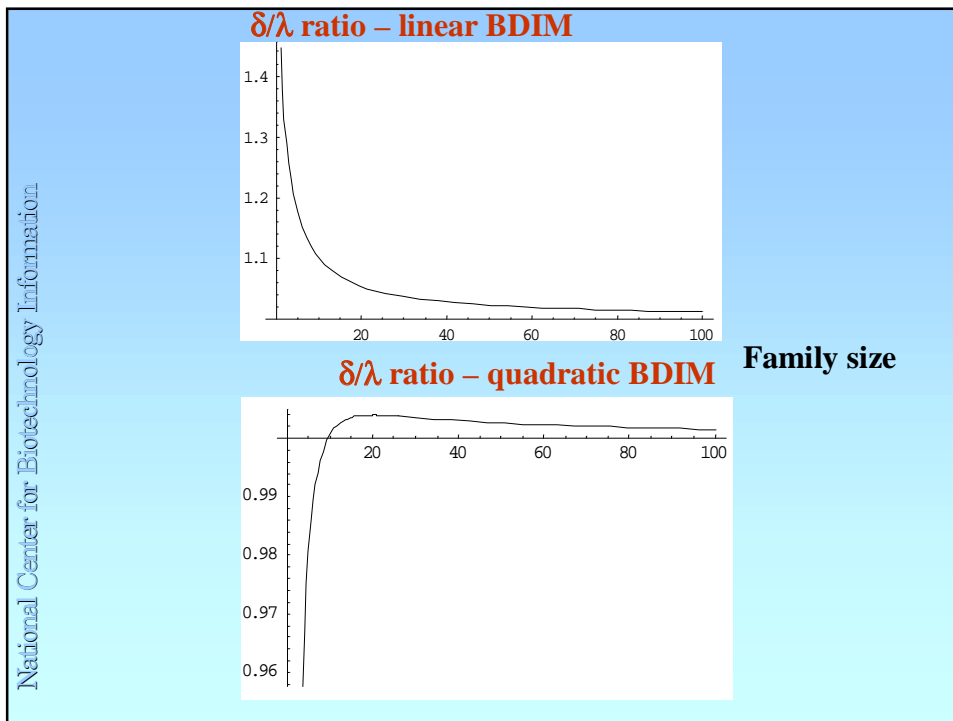
# Birth and Death of Domains: Simple Models of Genome Evolution Explain Power Law Distribution of Protein Family Size



Birth and Death of Domains: Simple Models of Genome Evolution Explain Power Law Distribution of Protein Family Size



# Birth and Death of Domains: Simple Models of Genome Evolution Explain Power Law Distribution of Protein Family Size



**Conclusions on stochastic (Markov) BDIMs and some general conclusions**

**2<sup>nd</sup> order balanced linear BDIM is sufficient to explain the observed distributions of domain family size**

**However, when a stochastic model is used to estimate the time required to reach the maximum family size, the estimate is ~10,000 times greater than the time suggested by empirical data**

**BDIMs of the degree 2-3 (quadratic/cubic) formally solve the problem by predicting evolutionary rate compatible with observations**

National Center for Biotechnology Information

**Conclusions on stochastic (Markov) BDIMs and some general conclusions**

**Biological interpretation of “interaction” between paralogs, which is intrinsic in higher-order BDIMs – does it reflect selection?**

**The dependence of birth and death rates on family size dramatically changes depending on BDIM order. This needs to be tested against detailed empirical analysis of paralogous families.**

National Center for Biotechnology Information



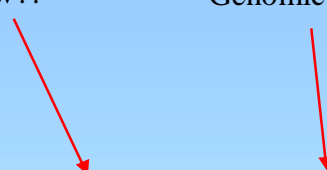
National Center for Biotechnology Information

**Some future directions...**

- **Incorporating selection into BDIM**
- **Combining BDIM with phylogenetic tree analysis**
- **Modeling evolutionary processes that lead to similar distributions in other contexts, e.g., multidomain protein architectures and interaction networks**

National Center for Biotechnology Information

Tomorrow??                      Genomics today



***“There are two kinds of science: physics and stamps collection”***

Attributed to Ernest Rutherford

Acknowledgments

National Center for Biotechnology Information

**Georgy Karev (NCBI, NIH),**  
**Yuri Wolf (NCBI, NIH),**  
Andrey Rzhetsky (Columbia University),  
Faina Berezovskaya (Howard University)