Entropy and diffusion characterize mutation accumulation and biological information loss

Stephan Baehr, 1*† Hans Baehr, 2,3,4*†

¹Biodesign Center for Mechanisms of Evolution, School of Life Sciences,
Arizona State University, USA

²Department of Physics and Astronomy,
University of Georgia, Athens, GA 30602, USA

³Center for Simulational Physics, University of Georgia, Athens, GA 30602, USA

⁴Max Planck Institute for Astronomy, Königstuhl 17, D-69117 Heidelberg, Germany

*E-mail: sbaehr@asu.edu, baehr@mpia.de † These authors contributed equally to this work.

Aging is a universal consequence of life, yet researchers have identified no universal theme. This manuscript considers aging from the perspective of entropy, wherein things fall apart. We first examine biological information change as a mutational distance, analogous to physical distance. In this model, informational change over time is fitted to an advection-diffusion equation, a normal distribution with a time component. The solution of the advection-diffusion equation provides a means of measuring the entropy of diverse biological systems. The binomial distribution is also sufficient to demonstrate that entropy increases as mutations or epimutations accumulate. As modeled, entropy scales with lifespans across the tree of life. This

perspective provides potential mechanistic insights and testable hypotheses as to how evolution has attained enhanced longevity: entropy management. We find entropy is an inclusive rather than exclusive aging theory.

Introduction

The biology of aging can be described by the phrase, "things fall apart". Researchers have noted that though there is rhyme and similarity to aging among individuals, each case is unique and unprogrammed [1, 2, 3]. Leonard Hayflick [4] has argued for decades that aging is entropy, an increase in molecular disorder over time. The concept thus far has only enjoyed modest popularity, perhaps because it does immediately offer a direct means of measurement, treatment, or a specific molecular mechanism.

References to entropy in aging research are often vague and allusory, lacking specific measurement and offering few testable hypotheses. In principle, the things that researchers have been measuring all along, and which are known to modify lifespan, should also fit naturally within the variables that describe entropy; this has yet to be emphatically shown. The purpose of this manuscript is to make a bridge between the world that the biology of aging knows and measures, and the physical understanding of entropy.

Superficially, the signature of entropy whittling at organism genomes across time may be recognized as the accumulation of deleterious mutations. As deleterious mutations accumulate, information disperses. Mutation accumulation is notably a familiar concept to aging biology, being among the oldest concepts in the field in one form or another. [5, 6, 7]. The idea of accumulating errors leading to critical breakdowns of biological systems is referred to as a "error catastrophe" or "mutation catastrophe". A modern version of mutation catastrophe is supported by some evidence; DNA mutations do accumulate over

time in cells over a lifespan.[8, 9, 10]. However, we note that mutation accumulation need not only refer to DNA, if the definition is broadened: what is epigenetic information loss, if it is not the accumulation of (epi)mutations? While acknowledging the importance of DNA mutations in the aging process [11], for example in the emergence of cancer [12], the field has been emphasizing for years that epimutations likely also have a proximal role to play in both aging and cancer. Compounding evidence [13] and recent experimentation [14] have highlighted the need for a model that includes, or even emphasizes the importance of epigenetic information in aging. Ideally, a theory should be flexible to account for aging of all sorts, even for organisms that age over the course of days, such as in *E. coli*.

We propose a model where the accumulation of mutations over time between at least two points can be considered a "mutational distance". We fit the concept of mutational distance to physics definitions of distance via an advection-diffusion equation for the Brownian motion of tracers in a fluid flow, and use the result to model the change in entropy over time. From a starting point of highly similar cells within a population, the cells accumulate mutational distance over time. The model fits to DNA mutation accumulation experiments. We then model epimutation as a primary factor in the determination of longevity, though the role of DNA mutation and any other system of entropic gain may be added to the model as appropriate. We fit the model to organisms of varying lifespan and demonstrate the model's flexibility, which predicts that an entropic failure threshold causes biological mortality, via age-related phenotypes. We also examine a simple binomial entropy conversion for diverse biological systems, and its application to age-related molecular change. These simplified models suggest that aging may be entropy; and that entropy also increases within germline lineages as well, in the relative absence of selection.

1 Results

The inspiration for this work begins by recognizing that the phenotypic outcomes of biological aging are shared with those of evolutionary biology's mutation accumulation (MA) experiments. In MA experiments, an increasing burden of random mutations results in phenotypic degradation, because the average DNA mutation is deleterious [15]. The sum of mutations per line is counted to provide an estimate for mutation rate. This mutation accumulation, or mutational distance, over evolutionary time and in the presence of selection leads to the differences that define individuals, populations, and species. At baseline, however, unchecked mutation accumulation leads to mutational meltdown, phenotypic degradation, and lineage extinction.

Within a population of cells, in an MA or within an aging soma, Figure 1 examines the behavior of mutations within a population of cells over time. When interpreting mutational distance as physical distance, this perspective allows the import of physics equations that model evolution over time, with particular consideration of Brownian motion, and random walks of molecules. The model considers mutation accumulation as a one-dimensional distance and how a population of molecules will be distributed as a function of time, following a normal distribution. The equation for a normally distributed variable x is:

$$f(x,\mu,\sigma) = \frac{1}{\sigma\sqrt{2\pi}}e^{-\frac{(x-\mu)^2}{2\sigma^2}} \tag{1}$$

where μ is the distribution mean and σ^2 is the variance.

An advection-diffusion equation is often used to model the diffusive spread of a quantity, for example a drop of food coloring in a stream, or a rubber duck race down a river, Enterrennen in German, considering the one-dimensional distance from a starting point. Figure 1A demonstrates the one-dimensional distribution of rubber duckies floating down a river with some current: the mean distance increases as a function of time, with some spread in the distribution, which follows a normal distribution. This solution to the advection-diffusion equation is simply a normal distribution with a time component t, diffusion coefficient D, current or flow rate $D\lambda$, and drag coefficient λ (see Appendix 2).

$$F(x,t) = \frac{1}{\sqrt{4\pi Dt}} e^{-(x-D\lambda t)^2/4Dt}.$$
 (2)

We note that as an MA experiment proceeds, the spread of the distribution widens as a function of the mean, which is conventionally characterized by a Poisson distribution. The widening of the distribution is well appreciated [16] and because the normal distribution is a good approximation of the Poisson distribution for a large enough Poisson mean, the normal distribution suffices. We test the advection-diffusion model upon real MA data arising from both WT and hypermutator $E.\ coli\ [16,\ 17]$ in Figure 1B. Figure 1B demonstrates that a reasonable goodness of fit is approximated by the solution to the advection-diffusion equation. The $D\lambda$ variable contains the fold-difference between the wild-type and hypermutator strains; about 110-fold, and may be considered analogous to the current or flow rate. The variable λ is a fitting parameter analogous to a drag coefficient, which helps fit the observed variance to the mean distance from a starting point of zero mutations.

In addition to DNA mutations, the term "epimutation" is now used to describe changes in informational content of the epigenome. In the same way that DNA mutation accumulation results in mutational distance, epimutation accumulation does the same. We consider the general proposal of the aging field, which is that epimutations within chromatin structure drive age-related phenotypes. Under this perspective, we consider the hypothesis that shorter lived organisms will have higher epimutation rates in aggregate, be they DNA methylation or histone marks, and longer-lived organisms will have lower

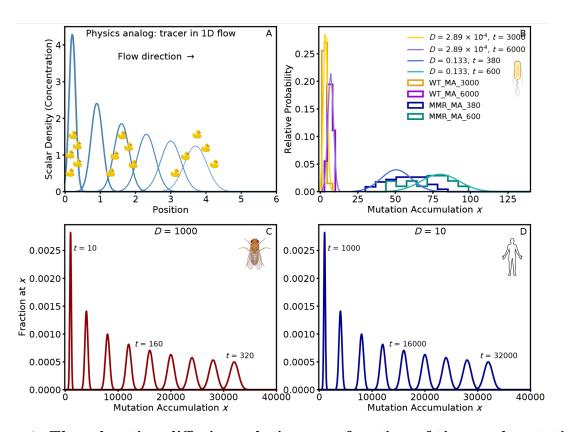


Figure 1: The advection-diffusion solution as a function of time and mutation accumulation. Panel A: A simple example of how a 1D diffusion model can model a passive tracer in a fluid, in this case, rubber duckies floating down a river. Panel B: The advection-diffusion equation applied to *E. coli* samples. Panels C and D: The advection-diffusion equation models how diverse organism lifespans may end up with the same final result. If chromatin "drift" or "epimutation" are primary drivers of aging, this equation is sufficiently flexible to model it.

epimutation rates in aggregate. Using epimutation burden estimates obtained from DNA methylation burden[18] as a coarse example, advection-diffusion model can be used to model epimutations under the same paradigm used to model DNA mutation rates. In Figure 1C and 1D, the model demonstrates that the same increase-in-variance outcome can be obtained over orders of magnitude of time-frame. The advection-diffusion equation is incredibly flexible, and therefore can be used to model informational distance change across the tree of life, and across biological informational storage media. We simply use a

rough DNA methylation average epimutation distance as a proof-of-principle for a broader epimutation rate argument.

Mapping and modeling biological information change over time is perhaps interesting, but a greater point can now be discussed with respect to entropy. A useful insight of the advection-diffusion equation is that it is Gaussian in nature, and from a Gaussian time series an estimate of the variance in terms of D and t can be achieved; which is independent of λ . This result, $\sigma^2 = 2Dt$, further expounded upon in the supplemental methods section, can be applied to an equation for the entropy of a Gaussian distribution, Equation (3). Therefore, technically any biological process with a normal distribution that is subject to change over time can be converted into units of entropy.

$$H = \frac{1}{2}(\ln(4\pi Dt) + 1). \tag{3}$$

Under a threshold model of system failure, such as has been frequently proposed in biology of aging research [7, 8], we hypothesize that organisms that reach their mortality sooner have higher rates of entropy gain, as measured by the variable D, over time; and long-lived organisms have lower rates of entropy gain within their systems. To this end, we extrapolate from our existing estimates of D and append short-lived E. coli[19] and long-lived Bristlecone pine trees $Pinus\ longaeva$ to demonstrate the effect of rate of variance spread upon estimates of entropy. The entropy equation models a log-linear relationship of increasing entropy with time. We note that the D extrapolated for E. coli in Figure 2 and the one directly calculated from DNA mutation accumulation data are highly divergent from one another; this may reflect the idea that $E\ coli$'s DNA mutation rate is sufficiently low, relative to its lifespan, that DNA mutation is incredibly unlikely to contribute to the replicative senescence experienced by the bacterium.

Entropy has been quantified in many biological systems and distributions [20]. Even

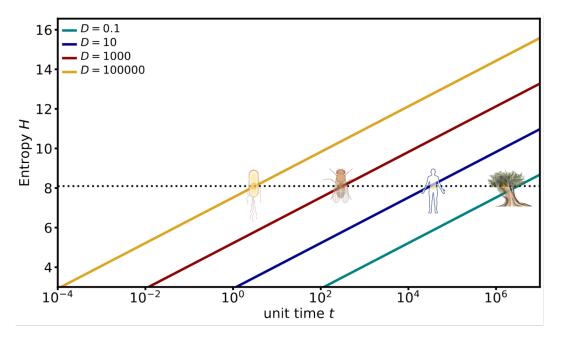


Figure 2: Comparison of different levels of diffusion on entropy as a function of time. With varying levels amounts of epigenetic diffusion (solid lines), the total entropy of a system will depend strongly on the value of the diffusion D. Orange points mark the rough maximum lifespans of a few organisms (from left to right): $E.\ coli,\ D.\ melanogaster,\ homo\ sapiens$ and $pinus\ longaeva$, assuming our unit of time is in days, as listed in Table 1.

the same system may be modeled by entropy in different different ways. For example, taking only a single cell of an aging tissue or from the *E. coli* MA, entropy can be quantified by the amount of information lost from acquiring 1, 2, or 80 mutations. Under the assumption that each mutation has an average effect and a simplified model assuming only two states, mutated or not, the entropy can be quantified by the binomial:

$$W = \binom{n}{k} = \frac{n!}{k! (n-k)!} = C(n,k)$$

$$\tag{4}$$

The solution of this equation further resolves into entropy via the unitless entropy equation (absent a Boltzmann constant)

$$S = \ln(W) \tag{5}$$

Where S is the entropy of a system in nats, and W is the number of microstates possible

from the accumulation of mutations in the *E. coli* genome. In the MA, the entropy becomes 15, 28, and 954 nats, respectively. It is unclear which measure of entropy, population level or within an individual cell, is most appropriate, or to what degree the entropy of a system needs to be scaled to make the measures equivalent. Regardless, in both cases calculated entropy increases, and both are directly linked to mutation rate. The general form of Figure 2 holds.

Instead of modeling only a single process, such as DNA mutation, there is a clear space in aging analysis for a parallelization of multiple parameters. The entropic gain of DNA mutation, transcript error, gene expression noise, translation error, chromatin structure, and even things like protein misfolding, or oxidative damage to proteins or lipids could in principle be modeled by some form of entropy. We summarize this perspective by slight modification of the advection diffusion equation from Equation (2) to incorporate all possible sources of entropy (i.e. DNA mutation, transcript error, translation error, chromatin structure, etc.)

$$H_{\text{total}} = \sum_{i} H_{i}. \tag{6}$$

where H_i is each potential source of entropy. The value of D and λ for each source of entropy can be independently calculated with distributions within a cell or set of cells, and we hypothesize that this quantity can be summed. At this time it is unclear as to how the varying system's entropic gains should be normalized; perhaps by some factor of the relative fitness cost per unit of entropy gained for each system measured. We simply note that the sum of all entropy within a cell or group of cells is relevant to encountering the entropic threshold, which we consider the "entropy catastrophe".

2 Discussion

The perspective of this manuscript examines a bridge between the biology of aging, physics, and evolutionary biology. This perspective began as a study of the phenotypic similarities emergent in mutation accumulation experiments and of the biology of aging, or specifically, an inescapable degradation of biological phenotypes in these contexts.

By recognizing that mutational distance is equivalent to physical distance, equations of physics may be applied to biological data. The advection-diffusion equation models the dispersion of molecules from a starting state of identical position, to a probability distribution of distance over time. This dispersion can be readily appreciated in the acquisition of DNA mutations over the course of human development and aging: all cells start out with an identical genotype, but over time a clock-like distance from the starting state to an aged state emerges over time. The model predicts an increased dispersion of mutational distance in chromatin or in DNA, though we also focus upon chromatin and its information storage role. Importantly, we find that the advection-diffusion equation can be modeled to fit mutation accumulation data from *E. coli* reasonably well, and can resolve the expected deviation in mutation rate between wild-type and hypermutator strains. The key insight of the advection-diffusion equation is its provision of a description of the variance of the Gaussian distribution in an age-related context.

Mathematicians and physicists have long recognized the importance of the Gaussian distribution to the study of entropy. The equation for entropy in a Gaussian system is straightforward (Equation (3)); the parameter of importance is the variable D in Equation (2). The model compares which values of D might give rise to known lifespan estimates across the tree of life. The model indicates that the organisms are hitting a similar entropy threshold for varying levels of D. By combining the perspectives of Figure 1

and Figure 2, a threshold model with a mean entropy acquisition rate over time, the model satisfyingly predicts even organisms with identical starting genotypes would reflect a Gaussian distribution of survivorship over time, focused about a mean. This is the result of aging experimentation on genetically identical organisms.[1]

The above perspective provides some avenues of application. The model proposed here, which we call the "entropy catastrophe hypothesis" for the biology of aging, provides testable hypotheses. Specifically, the hypothesis predicts that variance should increase over time, at least in the biological systems that are causal to aging. The accumulation of DNA mutations and epimutations are part of the hypothesis, but remain only two facets among the broad story of aging. The hypothesis predicts that evolutionary innovations that reduce entropy, such as increased replication fidelity, the induction of recycling programs, rewriting/restarting programs, and inducing purifying selection upon deleterious subsets of molecules, are responsible for enhanced longevity among organisms. The model proposes that interventions that increase the rate of aging, such as stress, temperature, or conditions like Hutchinson-Gilford progeria syndrome, ultimately act by increasing entropy in the system at biologically relevant levels.

There is undoubtedly a differential contribution to aging phenotypes and mortality, from differing molecular biological systems. Even within biological information systems, it likely true that the 'fitness impact' of epimutations in somatic cells is far less per mutation than that of DNA. Judging by relative mutation rates of the molecules as a proxy for relative importance, it may be that epimutations, 100 to 1000-fold more prevalent than DNA mutations, need to be weighted correspondingly such that each epimutation is 'worth' an inverse proportion to their prevalence; if not so extreme, it is certainly in that direction. To clarify, the reason translation or transcription errors are orders of magnitude more prevalent[21] than DNA is almost assuredly because their

individual impact is all that much less important than a single DNA mutation. The idea that chromatin information in mammals is the weak link is attractive to the field, but remains to be proven. For the present moment, we note that diffusion and entropy likely characterize the distributions of diverse age-related molecular phenomena. Their ultimate impacts, perhaps species and context specific, remain to be resolved.

Acknowledgments

The authors dedicate this work to their parents, Wolfgang Baehr and Jeanne Frederick. The authors thank their mentors, Michael Lynch and Hubert Klahr; and also their communities, scientific and otherwise, for their congenial atmosphere and *tolerance*. This work has been funded in part by NIH GMS grant 5R35GM122566-08 and the National Science Foundation, DBI-2119963, 2021-2026, BII: Mechanisms of Cellular Evolution. Figures 1 and 2 were created with the assistance of Biorender. Figure 2 was created with the assistance of Man Lin.

Methods and Supplementary information

Figure 1 has been generated in panel A by simple application of the advection-diffusion equation. Panel B experimental datasets are derived from the MA experiments reported in 2012[16] and 2025[17], whose data have been transformed/rearranged into Supplementary Table 2. Supplementary Table 1 provides the order-of-magnitude estimates of parameters that gave rise to figure 2; few *Drosophila* have ever measured to have a maximum lifespan of 100 days, but it is within a factor of 2 of reported values.

Panel C and D of 1 are estimates of epimutation rates derived from experimental results from 2023[18]; approximately 1% of DNA methylation sites are becoming discordant, or mutated, over several mammalian lifespans. For a similar amount in humans, 2.8×10^7 CpG sites results in 2.8×10^5 epimutations in a lifespan. If around 1-10% of those epimutations might be functional, the functional epimutational burden of DNA methylation alone may be on the order of 2,800 to 28,000. In contrast, the somatic DNA mutational burden of humans is on the order of 3,000-5,000, of which 1\%, or 30-50 mutations are functional. As an aside, the fitness effect-size and distribution of fitness effects for the epigenome are as yet unknown to our knowledge; but evolutionary theory predicts that they will be far less impactful, given their evolutionary impermanence. However, we may make an estimate of these based on their relative prevalence, and guess their average cellular fitness effect may be between 1/100th and 1/1000th that of the average DNA mutation. We acknowledge that the numbers offered are estimates, but offer the model as a general hypothesis to be tested. For simplicity of the model, this result has been translated into a relative number of human epimutations per aged cell; we suggest that the same principle and perhaps order of magnitude should extend also to chromatin marks, generally. Extending the epimutation hypothesis to D. melanogaster necessarily requires considering chromatin marks beyond DNA methylation, as flies lack DNA methylation.

The code used for the analysis of the data is available ¹.

Models

The following sections detail our mathematical methods and reasoning for assumptions of a random walk ansatz and by extension the advection-diffusion equation. From there we consider that a large number of persistent, dividing cells with accumulating chromatin epimutations obeys the central limit theorem, resulting in a Gaussian distribution of methylation states which can be modeled similar to a cloud of diffusing particles.

From this model, we formulate an expression for the total mutation load of a population in terms of the mutation rate, which depends on the measured dispersion of mutations across the population. We also evaluate an expression for the Shannon entropy of the information encoded within the nucleus of a single cell, which yields insights into the ways mutations accumulate within various organisms.

The Random Walk of a Single Genome

A random walk is a process by which something (i.e., a particle) can move from it's original location to a new position based on random, discrete movements. In one dimension, one can model the likelihood of displacement from a position using a binomial distribution B, where p is the probability of moving in the positive x direction, the additive inverse 1-p is movement in the negative x direction, k is the number of steps in the positive x direction and t is the number of trials or number of steps in the walk.

$$B(k, p, t) = \binom{k}{t} p^k (1 - p)^{t - k} \tag{7}$$

¹https://github.com/hbaehr/entropy

While this can be modeled in arbitrary dimensions we consider for now that a mutation in a sequence of base pairs or chromatin sites causes a cell within a population to 'walk' away from it's initial configuration in a single dimension x. At first, this allows for movement in the negative x direction, which is acceptable for modeling methylations, as the original configuration can be modified in one direction of more methylations but also in the opposite direction of fewer methylations. On the other hand, DNA mutations do not fit this framework as neatly, since the initial state of all base pairs can only 'move' in one direction: increasing mutation. However, we use this bidirectional random walk as an example that can be compared to a diffusive process and suggest limiting to only positive steps for the case of DNA mutations.

For enough trials (large t) or equivalently in this case, enough time, the binomial begins is well approximated by a Gaussian or normal distribution with mean $\mu = np$ and variance $\sigma^2 = np(1-p)$. Thus, while a binomial model works as a discrete distribution, also considering a continuous distribution allows us to draw a parallel to physical processes in fluid dynamics.

Epigenetic Evolution as a Diffusive Process

If we now look at a large collection of independently mutating epigenomes, such as a swath of skin cells, we can start to look at the large-scale pattern and evolution. We see a useful comparison with the evolution of a quantity that transports in one dimension through both advection and diffusion:

$$\frac{\partial}{\partial t}F(x,t) = D\left(\frac{\partial^2}{\partial x^2}F(x,t) + \lambda \frac{\partial}{\partial x}F(x,t)\right). \tag{8}$$

The first term on the right-hand side is the diffusive term where D is the diffusion constant and is assumed to be constant in time. Diffusion is a process that occurs when the net

motion of a group has some random component of the constituents. The second term is the advective term and $D\lambda$ is the advective (or drift) velocity, also constant in time and x. Advection has no random component and can be compared to a background flow field. We borrow the formulation of this equation which considers that drift in the positive x direction is due to a linear potential or a constant forcing $(\vec{F} = -\nabla U)$ which makes λ comparable to a drag or attenuation constant. In our case, it has the effect of adjusting the relative impact of diffusion or drift. For example, smaller values of λ will mean drift is less relevant to the displacement while higher values mean drift is more important. The solution to (8) is

$$F(x,t) = \frac{1}{\sqrt{4\pi Dt}} e^{-(x-D\lambda t)^2/4Dt},$$
(9)

where we use $\lambda = 0.1$ unless otherwise noted and caution against placing much physical significance into this value. The function F(x,t) represents the number of cells at time t within the population that have x number of changed methylation sites, centered around $x_0 = 0$ and $t_0 = 0$. Thus, we define this initial value problem by defining $F(0,0) = N\delta(x)$ where δ is the Dirac delta function and N is the number of cells in the population. This means that our final solution to the advective-diffusion equation is Equation (2) and means that the integral over the function is always N or in other words, all population members are represented somewhere along the distribution.

Mutations can arise from a number of sources, which we naively assume to be linear, such that the total population with mutations x at time t is

$$\sum_{i} F_i(x,t) = \sum_{i} \left(\frac{1}{\sqrt{4\pi D_i t}} e^{-(x-D_i \lambda_i t)^2/4D_i t} \right), \tag{10}$$

where the index i refers to different modes of mutation (i.e. chromatin, DNA, RNA, etc.).

Table 1: Model parameters

Model	genome size \mathcal{N}	diffusion D	λ	drift $(D\lambda)$	max. lifetime (days)
1 (pinus longaeva)	22×10^9	1/10	1/10	1/100	$4000 \text{ years} \times 365 = 1460000$
$2 (homo \ sapiens)$	3.2×10^9	10	1/10	1	$100 \text{ years} \times 365 = 36500$
3 (D. melanogaster)	180×10^{6}	1000	1/10	100	100
4 (e. coli)	4.6×10^{6}	100000	1/10	10000	3

Diffusion Coefficient

However, it would be useful to come up with a useful definition of D from laboratory data. We next seek to derive a value for the diffusion constant that makes sense for some model organism. We define D from Fick's law and the mean square displacement (MSD), which states that the displacement x from the initial position x_0 in one dimension at time t can be related to D as

$$\langle |x(t) - x_0|^2 \rangle = 2Dt, \tag{11}$$

where the angled brackets $\langle \cdot \rangle$ indicate an average over the entire population. However, since we have a uniform displacement this needs to be accounted for by subtracting $\langle x(t) \rangle^2$. This defines the mean distance between all the members of the group from their collective mean position, rather than their starting position:

$$\langle |x(t) - x_0|^2 \rangle - \langle |x(t)| \rangle^2 = 2Dt, \tag{12}$$

for the case of a diffusive model or equivalently for a binomial model

$$\langle |x(t) - x_0|^2 \rangle - \langle |x(t)| \rangle^2 = np(1-p). \tag{13}$$

From this definition we derive an approximation for the constant D with the data in Tables 2a through 2d of data for E. coli for two different strains at two different times. The first three come from [16] while the final dataset is measured in [17].

Table 2: Mutation accumulation counts across experiments.

Accumulation 0	1 2	3 4	5 6	7 8		Accu	mula	tion	4 5	6	7 8	9	10	11
Count 1	9 8	8 10	0 1	0 1		Cour	$_{ m t}$		2 1	6	2 5	2	2	1
(a) Wild type MA $t = 3000$					(b) Wild type MA $t = 6000$									
							,							
Accumulati	on 32	2 36	40	41	42	46	48	49	50	52	54	55	56	_
Count	1	1	2	1	1	1	2	1	1	1	1	3	1	
Accumulati	on 58	8 62	63	64	65	66	67	69	71	74	75	78	84	_
Count	1	1	1	3	1	2	1	1	1	1	1	1	1	
(c) MMR MA $t = 380$														
_	Accumulation Count		n 44	1 50	53	60	63	68	75	77				
			1	1	1	1	1	1	1	2				
_	A 1			70	00	01		0.0	0.0	1	_			
	Accumulation				80	81	85	86	96	155)			
	Count		0	1	1	1	1	1	1	1				
(d) MMR MA $t = 600$														

Shannon entropy

We now need a way of quantifying the information content or entropy, with the system. From information theory, the Shannon entropy H

$$H(X) \equiv -\sum_{x \in \mathcal{X}} p(x) \ln p(x), \tag{14}$$

describes the amount of uncertainty of the quality of information within the epigenome where p(x) is the probability or distribution of a state x [22].

We use both the binomial equation (7) and our solution to the advection-diffusion equation as a distribution of changes to epigenetic markers across a population of cells. For each cell, gene expression can be expressed as a distribution f that depends on a time interval t, the number of cell divisions N, and the natural variation of gene expression from one chromatin site to another.

The Shannon entropy of a binomially distributed random variable is

$$H_{\text{binom}} = \frac{1}{2} (\ln |np(1-p)2\pi| + 1). \tag{15}$$

while for a Gaussian distributed random variable it is

$$H_{\text{Gauss}} = \frac{1}{2} (\ln |\sigma^2 2\pi| + 1).$$
 (16)

One can see the similarity in the Shannon entropy for each model. For the Gaussian shapes introduced by the 'diffusion' of epigenetic mutations via the solution to the advectiondiffusion equation (Eq. (2)) where $\sigma^2 = 2Dt$, we arrive at an expression for the epigenetic entropy as defined in Equation (3). This assumes that D is constant in time and x, although there are many factors which could affect the value of D. An interesting feature of this formulation is that drift or advection $D\lambda$ only factors into the entropy gain through the diffusion constant with the factor λ omitted. To understand this we revisit the interpretation of λ . Our interpretation is that this represents a ratio of relative efficiency of drift versus diffusion and as such does not reflect on the nature of the system with information about either drift or diffusion. Furthermore, since entropy is the measure of the disorder of a system, λ contains no information about the distribution of states within the system. We can reconcile this by considering the simple case where D=0, which corresponds to the situation where all change occurs exactly on one methyl group (although not necessarily the same one) every unit of time in the same direction. As far as this model is concerned, the system of independently mutating cells retains its configuration for all times and thus has a constant entropy in time.

We plot the Shannon entropy for a few values of D in Figure 2 and compare with the approximate maximum lifespan of a few example organisms. We find that an entropy threshold of approximately 8 coincides with a number of these organisms. One possible interpretation is that, regardless of species, fitness breaks down at some entropy threshold

illustrated in Figure 2. What does change from one organism to another is the diffusion of epigenetic information, which can depend on a number of factors including but not limited to: epigenome size, body size, programming, repair mechanisms, and external (environmental) triggers. In Section ??, we speculate and explore possible ways to account for some of these factors in an advective-diffusive model of epigenetic evolution.

Accounting for Additional Mutagenic Effects

Our solution to the advection-diffusion equation permits various levels of flexibility to account for additional factors, such as a diffusion parameter that is not constant in time or space, source terms to account for external factors, etc. We therefore take a step back and look at the more general formation of the advection-diffusion equation

$$\frac{\partial}{\partial t}F(x,t) = D(t)\left(\frac{\partial^2}{\partial x^2}F(x,t) + \lambda \frac{\partial}{\partial x}F(x,t)\right) + \mathcal{S}(x,t),\tag{17}$$

where S is a source function that can represent the accumulation of epigenetic mutations from an external mechanism (for example, environmental factors such as radiation or carcinogen exposure) and D is now a function of t. Solutions with a non-zero source term can be found analytically, provided S has an exponential form similar to the solution (2). When D is a function of x, non-trivial solutions can only be found through numerical methods or also by parameterizing D(x) in terms of t.

One such example is When D is a function of time and S(x,t) = 0, the advectiondiffusion equation of (17) still has a fairly simple analytic solution. A time-variable diffusion could be used to explain declining repair mechanisms as an organism ages or the increase in mutagenicity of an organism with time. If one simply assumes a linearly increasing diffusion of (epi)genetic information, the Shannon entropy then increases quadratically in time, potentially drastically altering the increase in entropy as age in-

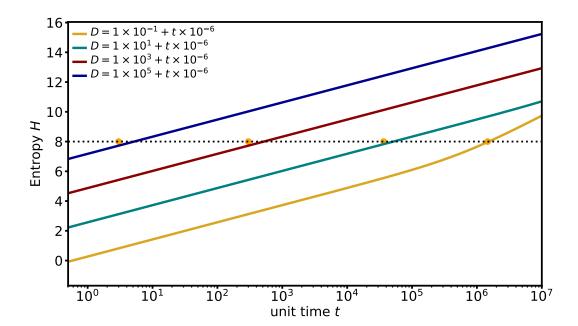


Figure 3: Comparison of different levels of linearly increasing diffusion on entropy. Same as Fig. 2, but diffusion increases linearly with time. The coefficient of the linear component is small such that only near the end of the least diffusive model is the increase in entropy noticeable.

creases as shown in Fig. 3.

References

- [1] C. Finch, T. B. L. Kirkwood, *Chance, development, and aging* (Oxford University Press, New York, 2000).
- [2] T. Kirkwood, S. Melov, Current Biology 21, R701 (2011).
- [3] U. Alon, Systems medicine: physiological circuits and the dynamics of disease, Computational biology series (CRC Press, Taylor & Francis Group, Boca Rato London New York, 2024), first edition edn.
- [4] L. Hayflick, *PLoS Genetics* **3**, e220 (2007).

- [5] P. B. Medawar, An Unsolved Problem of Biology (H. K. Lewis & Co. Ltd., London, 1952).
- [6] L. Szilard, Proceedings of the National Academy of Sciences 45, 30 (1959). Publisher: Proceedings of the National Academy of Sciences.
- [7] L. E. Orgel, Proceedings of the National Academy of Sciences 49, 517 (1963). Publisher: Proceedings of the National Academy of Sciences.
- [8] B. Milholland, Y. Suh, J. Vijg, Experimental Gerontology 94, 34 (2017). Publisher: Elsevier BV.
- [9] J. Vijg, Ageing Research Reviews **68**, 101316 (2021).
- [10] A. Cagan, et al., Nature **604**, 517 (2022).
- [11] C. López-Otín, M. A. Blasco, L. Partridge, M. Serrano, G. Kroemer, Cell 186, 243 (2023). Publisher: Elsevier BV.
- [12] C. Tomasetti, B. Vogelstein, *Science* **347**, 78 (2015). Publisher: American Association for the Advancement of Science (AAAS).
- [13] T. Rando, H. Chang, Cell 148, 46 (2012).
- [14] J.-H. Yang, et al., Cell 186, 305 (2023).
- [15] A. Eyre-Walker, P. D. Keightley, Nature Reviews. Genetics 8, 610 (2007).
- [16] H. Lee, E. Popodi, H. Tang, P. L. Foster, Proceedings of the National Academy of Sciences of the United States of America 109, E2774 (2012).
- [17] S. Baehr, et al., Genome Biology and Evolution 17, evaf049 (2025).

- [18] E. M. Bertucci-Richter, B. B. Parrott, Nature Communications 14, 7731 (2023).
- [19] P. Wang, et al., Current Biology 20, 1099 (2010). Publisher: Elsevier BV.
- [20] C. Adami, The Evolution of Biological Information: How Evolution Creates Complexity, from Viruses to Brains (Princeton University Press, 2024).
- [21] M. R. Lynch, Evolutionary cell biology (OUP, New York, 2023).
- [22] C. E. Shannon, Bell System Technical Journal 27, 379 (1948).