

Article

# The Free Energy Requirements of Biological Organisms; Implications for Evolution

David H. Wolpert <sup>1,2,3</sup>

<sup>1</sup> Santa Fe Institute, 1399 Hyde Park Road, Santa Fe, NM 87501, USA; david.h.wolpert@gmail.com

<sup>2</sup> Massachusetts Institute of Technology (MIT), 77 Massachusetts Ave, Cambridge, MA 02139, USA

<sup>3</sup> Arizona State University, Tempe, AZ 85281, USA

Academic Editors: John Baez, John Harte and Marc Harper

Received: 8 February 2016; Accepted: 8 April 2016; Published: 13 April 2016

**Abstract:** Recent advances in nonequilibrium statistical physics have provided unprecedented insight into the thermodynamics of dynamic processes. The author recently used these advances to extend Landauer’s semi-formal reasoning concerning the thermodynamics of bit erasure, to derive the minimal free energy required to implement an arbitrary computation. Here, I extend this analysis, deriving the minimal free energy required by an organism to run a given (stochastic) map  $\pi$  from its sensor inputs to its actuator outputs. I use this result to calculate the input-output map  $\pi$  of an organism that optimally trades off the free energy needed to run  $\pi$  with the phenotypic fitness that results from implementing  $\pi$ . I end with a general discussion of the limits imposed on the rate of the terrestrial biosphere’s information processing by the flux of sunlight on the Earth.

**Keywords:** thermodynamics of computation; Landauer bound; information processing rate of the biosphere

## 1. Introduction

It is a truism that biological systems acquire and store information about their environments [1–5]. However, they do not just store information; they also process that information. In other words, they perform computation. The energetic consequences for biological systems of these three processes—acquiring, storing, and processing information—are becoming the focus of an increasing body of research [6–15]. In this paper, I further this research by analyzing the energetic resources that an organism needs in order to compute in a fitness-maximizing way.

Ever since Landauer’s seminal work [16–26], it has been appreciated that the laws of statistical physics impose lower bounds on how much thermodynamic work must be done on a system in order for that system to undergo a two-to-one map, e.g., to undergo bit erasure. By conservation of energy, that work must ultimately be acquired from some external source (e.g., sunlight, carbohydrates, *etc.*). If that work on the system is eventually converted into heat that is dumped into an external heat bath, then the system acts as a heater. In the context of biology, this means that whenever a biological system (deterministically) undergoes a two-to-one map, it must use free energy from an outside source to do so and produces heat as a result.

These early analyses led to a widespread belief that there must be strictly positive lower bounds on how much free energy is required to implement any deterministic, logically-irreversible computation. Indeed, Landauer wrote “...logical irreversibility is associated with physical irreversibility and requires a minimal heat generation” [16]. In the context of biology, such bounds would translate to a lower limit on how much free energy a biological system must “harvest” from its environment in order to implement any particular (deterministic) computation, not just bit erasure.

A related conclusion of these early analyses was that a one-to-two map, in which noise is added to a system that is initially in one particular state with probability one, can act as a refrigerator, rather than a heater, removing heat from the environment [16,20–22]. Formally, the minimal work that needs to be done on a system in order to make it undergo a one-to-two map is negative. So for example, if the system is coupled to a battery that stores free energy, a one-to-two map can “power the battery”, by gaining free energy from a heat bath rather than dumping it there. To understand this intuitively, suppose we have a two-state system that is initially in one particular state with probability one. Therefore, the system initially has low entropy. That means we can connect it to a heat bath and then have it do work on a battery (assuming the battery was initially at less than maximum storage), thereby transferring energy from the heat bath into that battery. As it does this, though, the system gets thermalized, *i.e.*, undergoes a one-to-two map (as a concrete example, this is what happens in adiabatic demagnetization of an Ising spin system [16]).

This possibility of gaining free energy by adding noise to a computation, or at least reducing the amount of free energy the computation needs, means that there is a trade-off in biology: on the one hand, there is a benefit to having biological computation that is as precise as possible, in order to maximize the behavioral fitness that results from that computation; on the other hand, there is a benefit to having the computation be as imprecise as possible, in order to minimize the amount of free energy needed to implement that computation. This tradeoff raises the intriguing possibility that some biological systems have noisy dynamics “on purpose”, as a way to maintain high stores of free energy. For such a system, the noise would not be an unavoidable difficulty to be overcome, but rather a resource to be exploited.

More recently, there has been dramatic progress in our understanding of non-equilibrium statistical physics and its relation to information-processing [27–43]. Much of this recent literature has analyzed the minimal work required to drive a physical system’s (fine-grained) microstate dynamics during the interval from  $t = 0$  to  $t = 1$  in such a way that the associated dynamics over some space of (coarse-grained) macrostates is given by some specified Markov kernel  $\pi$ . In particular, there has been detailed analysis of the minimal work needed when there are only two macrostates,  $v = 0$  and  $v = 1$ , and we require that both get mapped by  $\pi$  to the macrostate  $v = 0$  [36,38,44]. By identifying the macrostates  $v \in V$  as Information Bearing Degrees of Freedom (IBDF) [22] of an information-processing device like a digital computer, these analyses can be seen as elaborations of the analyses of Landauer *et al.* on the thermodynamics of bit erasure. Recently, these analyses of maps over binary spaces  $V$  have been applied to explicitly biological systems, at least for the special case of a periodic forcing function [14].

These analyses have resulted in substantial clarifications of Landauer’s semiformal reasoning, arguably overturning it in some regards. For example, this analysis has shown that the logical (ir)reversibility of  $\pi$  has nothing to do with the thermodynamic (ir)reversibility of a system that implements  $\pi$ . In particular, it is possible to implement bit erasure (which is logically irreversible) in a thermodynamically-reversible manner. In the modern understanding, there is no irreversible increase of entropy in bit erasure. Instead, there is a minimal amount of thermodynamic work that needs to be expended in a (thermodynamically reversible) implementation of bit erasure (see Example 3 below.)

Many of these previous analyses consider processes for implementing  $\pi$  that are tailored for some specific input distribution over the macrostates,  $P(v_t)$ . Such processes are designed to be thermodynamically reversible when run on  $P(v_t)$ . However, when run on a distribution other than  $P(v_t)$ , they are thermodynamically irreversible, resulting in wasted (dissipated) work. For example, in [45], the amount of work required to implement  $\pi$  depends on an assumption for  $\epsilon$ , the probability of a one in a randomly-chosen position on the bit string.

In addition, important as they are, these recent analyses are not applicable to arbitrary maps  $\pi$  over a system’s macrostates. For example, as discussed in [46], the “quench-based” devices analyzed in [36,38,44] can only implement maps whose output is independent of its input (as an example, the output of bit erasure, an erased bit, is independent of the original state of the bit).

Similarly, the devices considered in [45,47] combine a “tape” containing a string of bits with a “tape head” that is positioned above one of the bits on the tape. In each iteration of the system, the bit currently under the tape head undergoes an arbitrary map to produce a new bit value, and then, the tape is advanced so that the system is above the next bit. Suppose that, inspired by [48], we identify the state of the IBDF of the overall tape-based system as the entire bit string, aligned so that the current tape position of the read/write subsystem is above Bit zero. In other words, we would identify each state of the IBDF as an aligned big string  $\{v_i : i = \dots, -1, 0, \dots, N\}$  where  $N$  is the number of bits that have already been processed, and the (negative) minimal index could either be finite or infinite (note that unless we specify which bit of the string is the current one, *i.e.*, which has index zero, the update map over the string is not defined).

This tape-based system is severely restricted in the set of computations it can implement on its IBDF. For example, because the tape can only move forward, the system cannot deterministically map an IBDF state  $v = \{\dots, v_{-1}, v_0, v_1, \dots, v_N\}$  to an IBDF state  $v' = \{\dots, v'_{-1}, v'_0, v'_1, \dots, v'_{N-1}\}$ . (In [49], the tape can rewind. However, such rewinding only arises due to thermal fluctuations and therefore does not overcome the problem.)

It should be possible to extend either the quench-based devices reviewed in [38] and the tape-based device introduced in [45] into a system that could perform arbitrary computation. In fact, in [46], I showed how to extend quench-based devices into systems that could perform arbitrary computation in a purely thermodynamically-reversible manner. This allowed me to calculate the minimal work that any system needs to implement any given conditional distribution  $\pi$ . To be precise, I showed how for any  $\pi$  and initial distribution  $P(v_t)$ , one could construct:

- a physical system  $\mathcal{S}$ ;
- a process  $\Lambda$  running over  $\mathcal{S}$ ;
- an associated coarse-grained set  $V$  giving the macrostates of  $\mathcal{S}$ ;

such that:

- running  $\Lambda$  on  $\mathcal{S}$  ensures that the distribution across  $V$  changes according to  $\pi$ , even if the initial distribution differs from  $P(v_t)$ ;
- $\Lambda$  is thermodynamically reversible if applied to  $P(v_t)$ .

By the second law, no process can implement  $\pi$  on  $P(v_t)$  with less work than  $\Lambda$  requires. Therefore, by calculating the amount of work required by  $\Lambda$ , we calculate a lower bound on how much work is required to run  $\pi$  on  $P(v_t)$ . In the context of biological systems, that bound is the minimal amount of free energy that any organism must extract from its external environment in order to run  $\pi$ .

However, just like in the systems considered previously in the literature, this  $\Lambda$  is thermodynamically optimized for that initial distribution  $P(v_t)$ . It would be thermodynamically irreversible (and therefore dissipate work) if used for any other other initial distribution. In the context of biological systems, this means that while natural selection may produce an information-processing organism that is thermodynamically optimal in one environment, it cannot produce one that is thermodynamically optimal in all environments.

Biological systems are not only information-processing systems, however. As mentioned above, they also acquire information from their environment and store it. Many of these processes have nonzero minimal thermodynamic costs, *i.e.*, the system must acquire some minimal free energy to implement them. In addition, biological systems often rearrange matter, thereby changing its entropy. Sometimes, these systems benefit by decreasing entropy, but sometimes, they benefit by increasing entropy, *e.g.*, as when cells use depletion forces, when they exploit osmotic pressures, *etc.* This is another contribution to their free energy requirements. Of course, biological systems also typically perform physical “labor”, *i.e.*, change the expected energy of various systems, by breaking/making chemical bonds, and on a larger scale, moving objects (including themselves),

developing, growing, *etc.* They must harvest free energy from their environment to power this labor, as well. Some biological processes even involve several of these phenomena simultaneously, e.g., a biochemical pathway that processes information from the environment, making and breaking chemical bonds as it does so and also changing its overall entropy.

In this paper, I analyze some of these contributions to the free energy requirements of biological systems and the implications of those costs for natural selection. The precise contributions of this paper are:

1. Motivated by the example of a digital computer, the analysis in [46] was formulated for systems that change the value  $v$  of a single set of physical variables,  $V$ . Therefore, for example, as formulated there, bit erasure means a map that sends both  $v_t = 0$  and  $v_t = 1$  to  $v_{t+1} = 0$ .

Here, I instead formulate the analysis for biological “input-output” systems that implement an arbitrary stochastic map taking one set of “input” physical variables  $X$ , representing the state of a sensor, to a separate set of “output” physical variables,  $Y$ , representing the action taken by the organism in response to its sensor reading. Therefore, as formulated in this paper, “bit erasure” means a map  $\pi$  that sends both  $x_t = 0$  and  $x_t = 1$  to  $y_{t+1} = 0$ . My first contribution is to show how to implement any given stochastic map  $X \rightarrow Y$  with a process that requires minimal work if it is applied to some specified distribution over  $X$  and to calculate that minimal work.

2. In light of the free energy costs associated with implementing a map  $\pi$ , what  $\pi$  would we expect to be favored by natural selection? In particular, recall that adding noise to a computation can result in a reduction in how much work is needed to implement it. Indeed, by using a sufficiently noisy  $\pi$ , an organism can increase its stored free energy (if it started in a state with less than maximal entropy). Therefore, noise might not just be a hindrance that an organism needs to circumvent; an organism may actually exploit noise, to “recharge its battery”. This implies that an organism will want to implement a “behavior”  $\pi$  that is noisy as possible.

In addition, not all terms in a map  $x_t \rightarrow y_{t+1}$  are equally important to an organism’s reproductive fitness. It will be important to be very precise in what output is produced for some inputs  $x_t$ , but for other inputs, precision is not so important. Indeed, for some inputs, it may not matter at all what output the organism produces in response. In light of this, natural selection would be expected to favor organisms that implement behaviors  $\pi$  that are as noisy as possible (thereby saving on the amount of free energy the organism needs to acquire from its environment to implement that behavior), while still being precise for those inputs where behavioral fitness requires it. I write down the equations for what  $\pi$  optimizes this tradeoff and show that it is approximated by a Boltzmann distribution over a sum of behavioral fitness and energy. I then use that Boltzmann distribution to calculate a lower bound on the maximal reproductive fitness over all possible behaviors  $\pi$ .

3. My last contribution is to use the preceding results to relate the free energy flux incident on the entire biosphere to the maximal “rate of computation” implemented by the biosphere. This relation gives an upper bound on the rate of computation that humanity as a whole can ever achieve, if it restricts itself to the surface of Earth.

In Section 2, I first review some of the basic quantities considered in nonequilibrium statistical physics and then review some of the relevant recent work in nonequilibrium statistical physics (involving “quenching processes”) related to the free energy cost of computation. I then discuss the limitations in what kind of computations that recent work can be used to analyze. I end by presenting an extension to that recent work that does not have these limitations (involving “guided quenching processes”). In Section 3, I use this extension to calculate the minimal free energy cost of any given input-output “organism”. I end this section by analyzing a toy model of the role that this free energy cost would play in natural selection. Those interested mainly in these biological implications can skip Section 2 and should still be able to follow the thrust of the analysis.

In this paper I extend the construction reviewed in [38] to show how to construct a system to perform any given computation in a thermodynamically reversible manner. (It seems likely that the tape-based system introduced in [45] could also be extended to do this.)

## 2. Formal Preliminaries

### 2.1. General Notation

I write  $|X|$  for the cardinality of any countable space  $X$ . I will write the Kronecker delta between any two elements  $x, x' \in X$  as  $\delta(x, x')$ . For any logical condition  $\zeta$ ,  $I(\zeta) = 1$  (0, respectively) if  $\zeta$  is true (false, respectively). When referring generically to any probability distribution, I will write “ $Pr$ ”. Given any distribution  $p$  defined over some space  $X$ , I write the Shannon entropy for countable  $X$ , measured in nats, as:

$$S_p(X) = - \sum_{x \in X} p(x) \ln [p(x)] \quad (1)$$

As shorthand, I sometimes write  $S_p(X)$  as  $S(p)$  or even just  $S(X)$  when  $p$  is implicit. I use similar notation for conditional entropy, joint entropy of more than one random variable, *etc.* I also write mutual information between two random variables  $X$  and  $Y$  in the usual way, as  $I(X; Y)$  [50–52].

Given a distribution  $q(x)$  and a conditional distribution  $\pi(x' | x)$ , I will use matrix notation to define the distribution  $\pi q$ :

$$[\pi q](x') = \sum_x \pi(x' | x) q(x) \quad (2)$$

For any function  $F(x)$  and distribution  $P(x)$ , I write:

$$\mathbb{E}_P(F) = \sum_x F(x) P(x) \quad (3)$$

I will also sometimes use capital letters to indicate variables that are marginalized over, e.g., writing:

$$\mathbb{E}_P(F(X, y)) = \sum_x P(x) F(x, y) \quad (4)$$

Below, I often refer to a process as “semi-static”. This means that these processes transform one Hamiltonian into another one so slowly that the associated distribution is always close to equilibrium, and as a result, only infinitesimal amounts of dissipation occur during the entire process. For this assumption to be valid, the implicit units of time in the analysis below must be sufficiently long on the timescale of the relaxation processes of the physical systems involved (or equivalently, those relaxation processes must be sufficiently quick when measured in those time units).

If a system with states  $x$  is subject to a Hamiltonian  $H(x)$ , then the associated equilibrium free energy is:

$$F_{eq}(H) \equiv -\beta^{-1} \ln[Z_H(\beta)] \quad (5)$$

where as usual  $\beta \equiv 1/kT$ , and the partition function is:

$$Z_H(\beta) = \sum_x \exp -\beta H(x) \quad (6)$$

However, the analysis below focuses on nonequilibrium distributions  $p(x)$ , for which the more directly relevant quantity is the nonequilibrium free energy, in which the distribution need not be a Boltzmann distribution for the current Hamiltonian:

$$\begin{aligned} F_{neq}(H, p) &\equiv \mathbb{E}_p(X) - kTS(p) \\ &= \sum_x p(x)H(x) + kT \sum_x p(x) \ln[p(x)] \end{aligned} \quad (7)$$

where  $k$  is Boltzmann's constant. For fixed  $H$  and  $T$ ,  $F_{neq}(H, p)$  is minimized by the associated Boltzmann distribution  $p$ , for which it has the value  $F_{eq}(H)$ . It will be useful below to consider the changes in nonequilibrium free energy that accompany a change from a distribution  $P$  to a distribution  $M$  accompanied by a change from a Hamiltonian  $H$  to a Hamiltonian  $H'$ :

$$\Delta F_{neq}^{H, H'}(P, M) \equiv F_{neq}(H', M) - F_{neq}(H, P) \quad (8)$$

## 2.2. Thermodynamically-Optimal Processes

If a process  $\Lambda$  maps a distribution  $P$  to a distribution  $M$  thermodynamically reversibly, then the amount of work it uses when applied to  $P$  is  $\Delta F_{neq}^{H, H'}(P, M)$  [38,48,53,54]. In particular,  $\Delta F_{neq}^{H, H'}(P, \pi P)$  is the amount of work used by a thermodynamically-reversible process  $\Lambda$  that maps a distribution  $P$  to  $\pi P$ . Equivalently, it is negative for the amount of work that is extracted by  $\Lambda$  when transforming  $P$  to  $\pi P$ .

In addition, by the second law, there is no process that maps  $P$  to  $M$  while requiring less work than a thermodynamically-reversible process that maps  $P$  to  $M$ . This motivates the following definition.

**Definition 1.** Suppose a system undergoes a process  $\Lambda$  that starts with Hamiltonian  $H$  and ends with Hamiltonian  $H'$ . Suppose as well that:

1. at both the start and finish of  $\Lambda$ , the system is in contact with a (single) heat bath at temperature  $T$ ;
2.  $\Lambda$  transforms any starting distribution  $P$  to an ending distribution  $\pi P$ , where neither of those two distributions need be at equilibrium for their respective Hamiltonians;
3.  $\Lambda$  is thermodynamically reversible when run on some particular starting distribution  $P$ .

Then,  $\Lambda$  is *thermodynamically optimal* for the tuple  $(P, \pi, H, H')$ .

**Example 1.** Suppose we run a process over a space  $X \times Y$ , transforming the  $t = 0$  distribution  $q(x)M(y)$  to a  $t = 1$  distribution  $p(x)M(y)$ . Therefore,  $x$  and  $y$  are statistically independent at both the beginning and the end of the process, and while the distribution over  $x$  undergoes a transition from  $q \rightarrow p$ , the distribution over  $y$  undergoes a cyclic process, taking  $M \rightarrow M$  (note that it is not assumed that the ending and starting  $y$ 's are the same or that  $x$  and  $y$  are independent at times between  $t = 0$  and  $t = 1$ ).

Suppose further that at both the beginning and end of the process, there is no interaction Hamiltonian, *i.e.*, at those two times:

$$H(x, y) = H^X(x) + H^Y(y) \quad (9)$$

Then, no matter how  $x$  and  $y$  are coupled during the process, no matter how smart the designer of the process, the process will require work of at least:

$$\Delta F_{neq}^{H, H'}(q, p) = \left( E_p(H^X) - E_q(H^X) \right) - kT \left( S(p) - S(q) \right) \quad (10)$$

Note that this amount of work is independent of  $M$ .

As a cautionary note, the work expended by any process operating on any initial distribution  $p(x)$  is the average of the work expended on each  $x$ . However, the associated change in nonequilibrium free energy is not the average of the change in nonequilibrium free energy for each  $x$ . This is illustrated in the following example.

**Example 2.** Suppose we have a process  $\Lambda$  that sends each initial  $x$  to an associated final distribution  $\pi(x' | x)$ , while transforming the initial Hamiltonian  $H$  into the final Hamiltonian  $H'$ . Write  $W_{H,H',\pi}^\Lambda(x)$  for the work expended by  $\Lambda$  when it operates on the initial state  $x$ . Then, the work expended by  $\Lambda$  operating on an initial distribution  $p(x)$  is  $\sum_x p(x)W_{H,H',\pi}^\Lambda(x)$ . In particular, choose the process  $\Lambda$ , so that it sends  $p \rightarrow \pi p$  with minimal work. Then:

$$\sum_x p(x)W_{H,H',\pi}^\Lambda(x) = \Delta F_{neq}^{H',H}(p, \pi p) \tag{11}$$

However, this does *not* equal the average over  $x$  of the associated changes to nonequilibrium free energy, *i.e.*,

$$\begin{aligned} \Delta F_{neq}^{H',H}(p, \pi p) &= F_{neq}(H', \pi p) - F_{neq}(H, p) \\ &\neq \sum_x p(x) \left[ F_{neq}(H', \pi(Y | x)) - F_{neq}(H, \delta(X, x)) \right] \end{aligned} \tag{12}$$

(where  $\delta(X, x)$  is the distribution over  $X$  that is a delta function at  $x$ ). The reason is that the entropy terms in those two nonequilibrium free energies are not linear; in general, for any probability distribution  $Pr(x)$ ,

$$\sum_x Pr(x) \ln[Pr(x)] \neq \sum_x Pr(x) \sum_{x'} \delta(x', x) \log[\delta(x', x)] \tag{13}$$

I now summarize what will be presented in the rest of this section.

Previous work showed how to construct a thermodynamically-optimal process for many tuples  $(p, \pi, H, H')$ . In particular, as discussed in the Introduction, it is known how to construct a thermodynamically-optimal process for any tuple  $(p, \pi, H, H')$  where  $\pi(x' | x)$  is independent of  $x$ , like bit erasure. Accordingly, we know the minimal work necessary to run any such tuple. In Section 2.3, I review this previous analysis and show how to apply it to the kinds of input-output systems considered in this paper.

However, as discussed in the Introduction, until recently, it was not known whether one could construct a thermodynamically-optimal process for any tuple  $(p, \pi, H, H')$ . In particular, given an arbitrary pair of an initial distribution  $p$  and conditional distribution  $\pi$ , it was not known whether there is a process  $\Lambda$  that is thermodynamically optimal for  $(p, \pi, H, H')$  for some  $H$  and  $H'$ . This means that it was not known what the minimal needed work is to apply an arbitrary stochastic map  $\pi$  to an arbitrary initial distribution  $p$ . In particular, it was not known if we could use the difference in nonequilibrium free energy between  $p$  and  $\pi p$  to calculate the minimal work needed to apply a computation  $\pi$  to an initial distribution  $p$ .

This shortcoming was overcome in [46], where it was explicitly shown how to construct a thermodynamically-optimal process for any tuple  $(p, \pi, H, H')$ . In Section 2.4, I show in detail how to construct such processes for any input-output system.

Section 2.4 also discusses the fact that a process that is thermodynamically optimal for  $(p, \pi, H, H')$  need not be thermodynamically optimal for  $(p', \pi, H, H')$  if  $p' \neq p$ . Intuitively, if we construct a process  $\Lambda$  that results in minimal required work for initial distribution  $p$  and conditional distribution  $\pi$ , but then apply that machine to a different distribution  $p' \neq p$ , then in general, work is dissipated. While that  $\Lambda$  is thermodynamically reversible when applied to  $p$ , in general, it is not

thermodynamically reversible when applied to  $p' \neq p$ . As an example, if we design a computer to be thermodynamically reversible for input distribution  $p$ , but then use it with a different distribution of inputs, then work is dissipated.

In a biological context, this means that if an organism is “designed” not to dissipate any work when it operates in an environment that produces inputs according to some  $p$ , but instead finds itself operating in an environment that produces inputs according to some  $p' \neq p$ , then it will dissipate extra work. That dissipated work is wasted since it does not change  $\pi$ , *i.e.*, has no consequences for the input-output map that the organism implements. However, by the conservation of energy, that dissipated work must still be acquired from some external source. This means that the organism will need to harvest free energy from its environment at a higher rate (to supply that dissipated work) than would an organism that were “designed” for  $p'$ .

### 2.3. Quenching Processes

A special kind of process, often used in the literature, can be used to transform any given initial nonequilibrium distribution into another given nonequilibrium distribution in a thermodynamically-reversible manner. These processes begin by quenching the Hamiltonian of a system. After that, the Hamiltonian is isothermally and quasi-statically changed, with the system in continual contact with a heat bath at a fixed temperature  $T$ . The process ends by applying a reverse quench to return to the original Hamiltonian (see [36,38,44] for discussion of these kinds of processes).

More precisely, such a *Quenching (Q) process* applied to a system with microstates  $r \in R$  is defined by:

1. an *initial/final* Hamiltonian  $H_{sys}^t(r)$ ;
2. an *initial* distribution  $\rho^t(r)$ ;
3. a *final* distribution  $\rho^{t+1}(r)$ ;

and involves the following three steps:

- (i) To begin, the system has Hamiltonian  $H_{sys}^t(r)$ , which is quenched into a first *quenching Hamiltonian*:

$$H_{quench}^t(r) \equiv -kT \ln[\rho^t(r)] \quad (14)$$

In other words, the Hamiltonian is changed from  $H_{sys}^t$  to  $H_{quench}^t$  too quickly for the distribution over  $r$  to change from  $\rho^t(r)$ .

Because the quench is effectively instantaneous, it is thermodynamically reversible and is adiabatic, involving no heat transfer between the system and the heat bath. On the other hand, while  $r$  is unchanged in a quench and, therefore, so is the distribution over  $R$ , in general, work is required if  $H_{quench}^t \neq H_{sys}^t$  (see [32,33,53,54]).

Note that if the Q process is applied to the distribution  $\rho^t$ , then at the end of this first step, the distribution is at thermodynamic equilibrium. However, if the process is applied to any other distribution, this will not be the case. In this situation, work is unavoidably dissipated in the next step.

- (ii) Next, we isothermally and quasi-statically transform  $H_{quench}^t$  to a second quenching Hamiltonian,

$$H_{quench}^{t+1}(r) \equiv -kT \ln[\rho^{t+1}(r)] \quad (15)$$

Physically, this means two things. First, that a smooth sequence of Hamiltonians, starting with  $H_{quench}^t$  and ending with  $H_{quench}^{t+1}$ , is applied to the system. Second, that while that sequence is being applied, the system is coupled with an external heat bath at temperature  $T$ , where the relaxation timescales of that coupling are arbitrarily small on the time scale of the dynamics of the

Hamiltonian. This second requirement ensures that to first order, the system is always in thermal equilibrium for the current Hamiltonian, assuming it started in equilibrium at the beginning of the step (recall from Section 2.1 that I assume that quasi-static transformations occur in an arbitrarily small amount of time, since the relaxation timescales are arbitrarily short).

- (iii) Next, we run a quench over  $R$  “in reverse”, instantaneously replacing the Hamiltonian  $H_{quench}^{t+1}(r)$  with the initial Hamiltonian  $H_{sys}^t$ , with no change to  $r$ . As in step (i), while work may be done (or extracted) in step (iii), no heat is transferred.

Note that we can specify any Q process in terms of its first and second quenching Hamiltonians rather than in terms of the initial and final distributions, since there is a bijection between those two pairs. This central role of the quenching Hamiltonians is the basis of the name “Q” process (I distinguish the distribution  $\rho$  that defines a Q process, which is instantiated in the physical structure of a real system, from the actual distribution  $P$  on which that physical system is run).

Both the first and third steps of any Q process are thermodynamically reversible, no matter what distribution that process is applied to. In addition, if the Q process is applied to  $\rho^t$ , the second step will be thermodynamically reversible. Therefore, as discussed in [36,38,48,54], if the Q process is applied to  $\rho^t$ , then the expected work expended by the process is given by the change in nonequilibrium free energy in going from  $\rho^t(r)$  to  $\rho^{t+1}(r)$ ,

$$\Delta F_{neq}^{H_{sys}^t, H_{sys}^t}(\rho^t, \rho^{t+1}) = \mathbb{E}_{\rho^{t+1}}(H_{sys}^t) - \mathbb{E}_{\rho^t}(H_{sys}^t) + kT \left[ S(\rho^t) - S(\rho^{t+1}) \right] \quad (16)$$

Note that because of how  $H_{quench}^t$  and  $H_{quench}^{t+1}$  are defined, there is no change in the nonequilibrium free energy during the second step of the Q process if it is applied to  $\rho^t$ :

$$\mathbb{E}_{\rho^{t+1}}(H_{quench}^{t+1}) - \mathbb{E}_{\rho^t}(H_{quench}^t) + kT \left[ S(\rho^t) - S(\rho^{t+1}) \right] = 0 \quad (17)$$

All of the work arises in the first and third steps, involving the two quenches.

The relation between Q processes and information-processing of macrostates arises once we specify a partition over  $R$ . I end this subsection with the following example of a Q process:

**Example 3.** Suppose that  $R$  is partitioned into two bins, *i.e.*, there are two macrostates. For both  $t = 0$  and  $t = 1$ , for both partition elements  $v$ , with abuse of notation, define:

$$P^t(v) \equiv \sum_{r \in v} \rho^t(r | v) \quad (18)$$

so that:

$$\rho^t(r) = \sum_v P^t(v) \rho^t(r | v) \quad (19)$$

Consider the case where  $P^0(v)$  has full support, but  $P^1(v) = \delta(v, 0)$ . Therefore, the dynamics over the macrostates (bins) from  $t = 0$  to  $t = 1$  sends both  $v$ 's to zero. In other words, it erases a bit.

For pedagogical simplicity, take  $H_{sys}^0 = H_{sys}^1$  to be uniform. Then, plugging in to Equation (16), we see that the minimal work is:

$$\begin{aligned}
 kT[S(\rho^0) - S(\rho^1)] &= kT \left[ S(P^0) + \sum_v P^0(v) \left( - \sum_r P^0(r | v) \ln[\rho(r | v)] \right) \right] - \{0 \rightarrow 1\} \\
 &= kT \left[ S(P^0) + \sum_{v^0} P^0(v) S(R^0 | v^0) \right] - \{0 \rightarrow 1\} \\
 &= kT \left[ S(P^0) + S(R^0 | V^0) - S(P^1) - S(R^1 | V^1) \right] \\
 &= kT \left[ S(P^0) + S(R^0 | V^0) - S(R^1 | V^1) \right] \tag{20}
 \end{aligned}$$

(the two terms  $S(R^t | v^t)$  are sometimes called “internal entropies” in the literature [38]).

In the special case that  $P^0(v)$  is uniform and that  $S(R^t | v^t)$  is the same for both  $t$  and both  $v_t$ , we recover Landauer’s bound,  $kT \ln(2)$ , as the minimal amount of work needed to erase the bit. Note though that outside of that special case, Landauer’s bound does not give the minimal amount of work needed to erase a bit. Moreover, in all cases, the limit in Equation (20) is on the amount of work needed to erase the bit; a bit can be erased with zero dissipated work, *pace* Landauer. For this reason, the bound in Equation (20) is sometimes called “generalized Landauer cost” in the literature [38].

On the other hand, suppose that we build a device to implement a Q process that achieves the bound in Equation (20) for one particular initial distribution over the value of the bit,  $\mathcal{G}_0(v)$ . Therefore, in particular, that device has “built into it” a first and second quenching Hamiltonian given by:

$$H_{quench}^0(r) = -kT \ln[\mathcal{G}_0(r)] \tag{21}$$

$$H_{quench}^1(r) = -kT \ln[\mathcal{G}_1(r)] \tag{22}$$

respectively, where:

$$\mathcal{G}_0(r) \equiv \sum_v \mathcal{G}_0(v) \rho^0(r | v) \tag{23}$$

$$\mathcal{G}_1(r) \equiv \rho^1(r | v = 0) \tag{24}$$

If we then apply that device with a different initial macrostate distribution,  $\mathcal{P}_1(v) \neq \mathcal{G}_0(v)$ , in general, work will be dissipated in step (ii) of the Q process, because  $\mathcal{P}_1(r) = \sum_v \mathcal{P}_1(v) \rho^0(r | v)$  will not be an equilibrium for  $H_{quench}^0$ . In the context of biology, if a bit-erasing organism is optimized for one environment, but then used in a different one, it will necessarily be inefficient, dissipating work (the minimal amount of work dissipated is given by the drop in the value of the Kullback–Leibler divergence between  $\mathcal{G}_t$  and  $\mathcal{P}_t$  as the system develops from  $t = 0$  to  $t = 1$ ; see [46]).

#### 2.4. Guided Q Processes

Soon after the quasi-static transformation step of any Q process begins, the system is thermally relaxed. Therefore, all information about  $r_t$ , the initial value of the system’s microstate, is quickly removed from the distribution over  $r$  (phrased differently, that information has been transferred into inaccessible degrees of freedom in the external heat bath). This means that the second quenching Hamiltonian cannot depend on the initial value of the system’s microstate; after that thermal relaxation of the system’s microstate, there is no degree of freedom in the microstate that has any information concerning the initial microstate. This means that after the relaxation, there is no degree of freedom within the system undergoing the Q process that can modify the second quenching Hamiltonian based on the value of the initial microstate.

As a result, *by itself*, a Q process cannot change an initial distribution in a way that depends on that initial distribution. In particular, it cannot map different initial macrostates to different

final macrostates (formally, a Q process cannot map a distribution with support restricted to the microstates in the macrostate  $v_t$  to one final distribution and map a distribution with support restricted to the macrostate  $v'_t \neq v_t$  to a different final distribution).

On the other hand, both quenching Hamiltonians of a Q process running on a system  $\mathcal{R}$  with microstates  $r \in R$  can depend on  $s_t \in S$ , the initial microstate of a different system,  $\mathcal{S}$ . Loosely speaking, we can run a process over the joint system  $\mathcal{R} \times \mathcal{S}$  that is thermodynamically reversible and whose effect is to implement a different Q process over  $R$ , depending on the value  $s_t$ . In particular, we can “coarse-grain” such dependence on  $s_t$ : given any partition over  $S$  whose elements are labeled by  $v \in V$ , it is possible that both quenching Hamiltonians of a Q process running on  $\mathcal{R}$  are determined by the macrostate  $v_t$ .

More precisely, a *Guided Quenching (GQ) process* over  $R$  guided by  $V$  (for conditional distribution  $\bar{\pi}$  and initial distribution  $\rho^t(r, s)$ )” is defined by a quadruple:

1. an *initial/final* Hamiltonian  $H^t_{sys}(r, s)$ ;
2. an *initial* joint distribution  $\rho^t(r, s)$ ;
3. a time-independent partition of  $S$  specifying an associated set of macrostates,  $v \in V$ ;
4. a conditional distribution  $\bar{\pi}(r | v)$ .

It is assumed that for any  $s, s'$  where  $s \in V(s')$ ,

$$\rho^t(r | s) = \rho^t(r | s') \tag{25}$$

*i.e.*, that the distribution over  $r$  at the initial time  $t$  can depend on the macrostate  $v$ , but not on the specific microstate  $s$  within the macrostate  $v$ . It is also assumed that there are boundary points in  $S$  (“potential barriers”) separating the members of  $V$  in that the system cannot physically move from  $v$  to  $v' \neq v$  without going through such a boundary point.

The associated GQ process involves the following steps:

- (i) To begin, the system has Hamiltonian  $H^t_{sys}(r, s)$ , which is quenched into a first quenching Hamiltonian written as:

$$H^t_{quench}(r, s) \equiv H^t_{quench;S}(s) + H^t_{quench;int}(r, s) \tag{26}$$

We take:

$$H^t_{quench;int}(r, s) \equiv -kT \ln[\rho^t(r | s)] \tag{27}$$

and for all  $s$  except those at the boundaries of the partition elements defining the macrostates  $V$ ,

$$H^t_{quench;S}(s) \equiv -kT \ln[\rho^t(s)] \tag{28}$$

However, at the  $s$  lying on the boundaries of the partition elements defining  $V$ ,  $H^t_{quench;S}(s)$  is arbitrarily large. Therefore, there are infinite potential barriers separating the macrostates of  $\mathcal{S}$ .

Note that away from those boundaries of the partition elements defining  $V$ ,  $\rho^t(r, s)$  is the equilibrium distribution for  $H^t_{quench}$ .

- (ii) Next, we isothermally and quasi-statically transform  $H^t_{quench}$  to a second quenching Hamiltonian,

$$H^{t+1}_{quench;S}(r, s) \equiv H^t_{quench;S}(s) + H^{t+1}_{quench;int}(r, s) \tag{29}$$

where:

$$H^{t+1}_{quench;int}(r, s) \equiv -kT \ln[\bar{\pi}(r | V(s))] \tag{30}$$

( $V(s)$  being the partition element that contains  $s$ ).

Note that the term in the Hamiltonian that only concerns  $\mathcal{S}$  does not change in this step. Therefore, the infinite potential barriers delineating partition boundaries in  $S$  remain for the entire step. I assume that as a result of those barriers, the coupling of  $\mathcal{S}$  with the heat bath during this step cannot change the value of  $v$ . As a result, even though the distribution over  $r$  changes in this step, there is no change to the value of  $v$ . To describe this, I say that  $v$  is “semi-stable” during this step. (To state this assumption more formally, let  $A(s', s'')$  be the (matrix) kernel that specifies the rate at which  $s' \rightarrow s''$  due to heat transfer between  $\mathcal{S}$  and the heat bath during this step (ii) [32,33]. Then, I assume that  $A(s', s'')$  is arbitrarily small if  $V(s'') \neq V(s')$ .)

As an example, the different bit strings that can be stored in a flash drive all have the same expected energy, but the energy barriers separating them ensure that the distribution over bit strings relaxes to the uniform distribution infinitesimally slowly. Therefore, the value of the bit string is semi-stable.

Note that even though a semi-stable system is not at thermodynamic equilibrium during its “dynamics” (in which its macrostate does not change), that dynamics is thermodynamically reversible, in that we can run it backwards in time without requiring any work or resulting in heat dissipation.

- (iii) Next, we run a quench over  $R \times S$  “in reverse”, instantaneously replacing the Hamiltonian  $H_{quench}^{t+1}(r, s)$  with the initial Hamiltonian  $H_{sys}^t(r, s)$ , with no change to  $r$  or  $s$ . As in step (i), while work may be done (or extracted) in step (iii), no heat is transferred.

There are two crucial features of GQ processes. The first is that a GQ process faithfully implements  $\bar{\pi}$  even if its output varies with its input and does so no matter what the initial distribution over  $R \times S$  is. The second is that for a particular initial distribution over  $R \times S$ , implicitly specified by  $H_{quench}^t(r, s)$ , the GQ process is thermodynamically reversible.

The first of these features is formalized with the following result, proven in Appendix A:

**Proposition 1.** A GQ process over  $R$  guided by  $V$  (for conditional distribution  $\bar{\pi}$  and initial distribution  $\rho^t(r, s)$ ) will transform any initial distribution  $p^t(v)p^t(r | v)$  into a distribution  $p^t(v)\bar{\pi}(r | v)$  without changing the distribution over  $s$  conditioned on  $v$ .

Consider the special case where the GQ process is in fact applied to the initial distribution that defines it,

$$\rho^t(r, s) = \sum_v \rho^t(v)\rho^t(s | v)\rho^t(r | v) \tag{31}$$

(recall Equation (25)). In this case, the initial distribution is a Boltzmann distribution for the first quenching Hamiltonian; the final distribution is:

$$\rho^{t+1}(r, s) = \sum_v \rho^t(v)\rho^t(s | v)\bar{\pi}(r | v) \tag{32}$$

and the entire GQ process is thermodynamically reversible. This establishes the second crucial feature of GQ processes.

Plugging in, in this special case, the change in nonequilibrium free energy is:

$$\Delta F_{neq}^{H_{sys}^t, H_{sys}^{t+1}}(\rho^t, \rho^{t+1}) = \left[ \sum_{r,s,v} \rho^t(v)\rho^t(s | v)(\bar{\pi}(r | v) - \rho^t(r | v))H_{sys}^t(r, s) \right] - kT \left[ S(\rho^{t+1}) - S(\rho^t) \right] \tag{33}$$

This is the minimal amount of free energy needed to implement the GQ process. An important example of such a thermodynamically-optimal GQ process is the work-free copy process discussed in [38] and the references therein.

Suppose that we build a device to implement a GQ process over  $R$  guided by  $V$  for conditional distribution  $\bar{\pi}$  and initial distribution:

$$\rho^t(r, s) = \sum_v \rho^t(r | v) \rho^t(s | v) \mathcal{G}_t(v) \quad (34)$$

Therefore, that device has “built into it” first and second quenching Hamiltonians that depend on  $\rho^t(r | v)$ ,  $\rho^t(s | v)$  and  $\mathcal{G}_t$ . Suppose we apply that device in a situation where the initial distribution over  $r$  conditioned on  $v$  is in fact  $\rho^t(r | v)$  and the initial distribution over  $s$  conditioned on  $v$  is in fact  $\rho^t(s | v)$ , but the initial macrostate distribution,  $P_t(v)$ , does not equal  $\mathcal{G}_t(v)$ . In this situation, the actual initial distribution at the start of step (ii) of the GQ process will not be an equilibrium for the initial quenching Hamiltonian. However, this will not result in there being any work dissipated during the thermal relaxation of that step. That is because the distribution over  $v$  in that step does not relax, no matter what it is initially (due to the infinite potential barriers in  $S$ ), while the initial distribution over  $(r, s)$  conditioned on  $v$  is in thermal equilibrium for the initial quenching Hamiltonian.

However, now suppose that we apply the device in a situation where the initial distribution over  $r$  conditioned on  $v$  does not equal  $\rho^t(r | v)$ . In this situation, work will be dissipated in step (ii) of the GQ process. That is because the initial distribution over  $r$  when the relaxation starts is not in thermal equilibrium for the initial quenching Hamiltonian, and this distribution does relax in step (ii). Therefore, if the device was not “designed” for the actual initial distribution over  $r$  conditioned on  $v$  (*i.e.*, does not use a  $\rho^t(r | v)$  that equals that actual distribution), it will necessarily dissipate work.

As elaborated below, this means that if a biological organism that implements any map  $\bar{\pi}$  is optimized for one environment, *i.e.*, one distribution over its inputs, but then used in an environment with a different distribution over its inputs, it will necessarily be inefficient, dissipating work (recall that above, we established a similar result for the specific type of Q process that can be used to erase a bit).

### 3. Organisms

In this section, I consider biological systems that process an input into an output, an output that specifies some action that is then taken back to the environment. As shorthand, I will refer to any biological system that does this as an “organism”. A cell exhibiting chemotaxis is an example of an organism, with its input being (sensor readings of) chemical concentrations and its output being chemical signals that in turn specify some directed motion it will follow. Another example is a eusocial insect colony, with its inputs being the many different materials that are brought into the nest (including atmospheric gases) and its output being material waste products (including heat) that in turn get transported out of the colony.

Physically, each organism contains an “input subsystem”, a “processor subsystem” and an “output subsystem” (among others). The initial macrostate of the input subsystem is formed by sampling some distribution specified by the environment and is then copied to the macrostate of the processor subsystem. Next, the processor iterates some specified first-order time-homogenous Markov chain (for example, if the organism is a cell, this Markov chain models the iterative biochemical processing of the input that takes place within the organism). The ending value of the chain is the organism’s output, which specifies the action that the organism then takes back to its environment. In general, it could be that for certain inputs, an organism never takes any action back to its environment, but instead keeps processing the input indefinitely. Here, that is captured by having the Markov chain keep iterating (*e.g.*, the biochemical processing keeps going) until it produces a value that falls within a certain predefined *halting* (sub)set, which is then copied to the organism’s output (the possibility that the processing never halts also ensures that the organism is Turing complete [55–57]).

There are many features of information processing in real biological systems that are distorted in this model; it is just a starting point. Indeed, some features are absent entirely. In particular, since

the processing is modeled as a first-order Markov chain, there is no way for an organism described by this model to “remember” a previous input it received when determining what action to take in response to a current input. Such features could be incorporated into the model in a straight-forward way and are the subject of future work.

In the next subsection, I formalize this model of a biological input-output system, in terms of an input distribution, a Markov transition matrix and a halting set. I then analyze the minimal amount of work needed by any physical system that implements a given transition matrix when receiving inputs from a given distribution, *i.e.*, the minimal amount of work a real organism would need to implement its input-output behavior that it exhibits in its environment, if it were free to use any physical process that obeys the laws of physics. To perform this analysis, I will construct a specific physical process that implements an iteration of the Markov transition matrix of a given organism with minimal work, when inputs are generated according to the associated input distribution. This process involves a sequence of multiple GQ processes. *It cannot be emphasized enough that these processes I construct are not intended to describe what happens in real biological input-output systems, even as a cartoon.* These processes are used only as a calculational tool, for finding a lower bound on the amount of work needed by a real biological organism to implement a given input-output transition matrix.

Indeed, because real biological systems are often quite inefficient, in practice, they will often use far more work than is given by the bound I calculate. However, we might expect that in many situations, the work expended by a real biological system that behaves according to some transition matrix is approximately proportional to the work that would be expended by a perfectly efficient system obeying the same transition matrix. Under that approximation, the relative sizes of the bounds given below should reflect the relative sizes of the amounts of work expended by real biological systems.

### 3.1. The Input and Output Spaces of an Organism

Recall from Section 2.4 that a subsystem  $S$  cannot use a thermodynamically-reversible Q process to update its own macrostate in an arbitrary way. However a different subsystem  $S'$  can guide an arbitrary updating of the macrostate of  $S$ , with a GQ process. In addition, the work required by a thermodynamically-reversible process that implements a given conditional distribution from inputs to outputs is the same as the work required by any other thermodynamically-reversible process that implements that same distribution.

In light of these two facts, for simplicity, I will not try to construct a thermodynamically-reversible process that implements any given organism’s input-output distribution directly, by iteratively updating the processor until its state lies in the halting subset and then copying that state to the output. Instead, I will construct a thermodynamically-reversible process that implements that same input-output distribution, but by “ping-ponging” GQ processes back and forth between the state of the processor and the state of the output system, until the output’s state lies in the halting set.

Let  $W$  be the space of all possible microstates of a *processor* subsystem, and  $U$  the (disjoint) space of all possible microstates of an *output* subsystem. Let  $\mathcal{X}$  be a partition of  $W$ , *i.e.*, a coarse-graining of it into a countable set of macrostates. Let  $X$  be the set of labels of those partition elements, *i.e.*, the range of the map  $\mathcal{X}$  (for example, in a digital computer,  $\mathcal{X}$  could be a map taking each microstate of the computer’s main RAM,  $w \in W$ , into the associated bit string,  $\mathcal{X}(w) \in X$ ). Similarly, let  $\mathcal{Y}$  be a partition of  $U$ , the microstate of the output subsystem. Let  $Y$  be the set of labels of those partition elements, *i.e.*, the range of the map  $\mathcal{Y}$ , with  $Y_{halt} \subseteq Y$  the halting subset of  $Y$ . I generically write an element of  $X$  as  $x$  and an element of  $Y$  as  $y$ . I assume that  $X$  and  $Y$ , the spaces of labels of the processor and output partition elements, respectively, have the same cardinality and, so, indicate their elements with the same labels. In particular, if we are concerned with Turing-complete organisms,  $X$  and  $Y$  would both be  $\{0, 1\}^*$ , the set of all finite bit strings (a set that is bijective with  $\mathbb{N}$ ).

For notational convenience, I arbitrarily choose one non-empty element of  $X$  and one non-empty element of  $Y$  and the additional label 0 to both of them (for example, in a Turing machine, it could be that we assign the label 0 to the partition element that also has label  $\{0\}$ ). Intuitively, these elements represent the “initialized” state of the processor and output subsystems, respectively.

The biological system also contains an *input* subsystem, with microstates  $f \in F$  and coarse-graining partition  $\mathcal{F}$  that produces macrostates  $b \in B$ . The space  $B$  is the same as the space  $X$  (and therefore is the same as  $Y$ ). The state of the input at time  $t = 0$ ,  $b_0$ , is formed by sampling an *environment distribution*  $\mathcal{P}_1$ . As an example,  $b_0$  could be determined by a (possibly noisy) sensor reading of the external environment. As another example, the environment of an organism could directly perturb the organism’s input macrostate at  $t = 0$ . For simplicity, I assume that both the processor subsystem and the output subsystem are initialized before  $b_0$  is generated, *i.e.*, that  $x_0 = y_0 = 0$ .

After  $b_0$  is set this way, it is copied to the processor subsystem, setting  $x_1$ . At this point, we iterate a sequence of GQ processes in which  $x$  is mapped to  $y$ , then  $y$  is mapped to  $x$ , then that new  $x$  is mapped to a new  $y$ , *etc.*, until (and if)  $y \in Y_{halt}$ . To make this precise, adopt the notation that  $[\alpha, \alpha']$  refers to the joint state ( $x = \alpha, y = \alpha'$ ). Then, after  $x_1$  is set, we iterate the following multi-stage *ping-pong* sequence:

1.  $[x_t, 0] \rightarrow [x_t, y_t]$ , where  $y_t$  is formed by sampling  $\pi(y_t | x_t)$ ;
2.  $[x_t, y_t] \rightarrow [0, y_t]$ ;
3. If  $y_t \in Y_{halt}$ , the process ends;
4.  $[0, y_t] \rightarrow [y_t, y_t]$ ;
5.  $[y_t, y_t] \rightarrow [y_t, 0]$ ;
6. Return to (1) with  $t$  replaced by  $t + 1$ ;

If this process ends (at stage (3)) with  $t = \tau$ , then the associated value  $y_\tau$  is used to specify an action by the organism back on its environment. At this point, to complete a thermodynamic cycle, both  $x$  and  $y$  are reinitialized to zero, in preparation for a new input.

Here, for simplicity, I do not consider the thermodynamics of the physical system that sets the initial value of  $b_0$  by “sensing the environment”; nor do I consider the thermodynamics of the physical system that copies that value to  $x_0$  (see [38] and the references therein for some discussion of the thermodynamics of copying). In addition I do not analyze the thermodynamics of the process in which the organism uses  $y_\tau$  to “take an action back to its environment” and thereby reinitializes  $y$ . I only calculate the minimal work required to implement the phenotype of the organism, which here is taken to mean the iterated ping-pong sequence between  $X$  and  $Y$ .

Moreover, I do not make any assumption for what happens to  $b_0$  after it is used to set  $x_1$ ; it may stay the same, may slowly decay in some way, *etc.* Accordingly, none of the thermodynamic processes considered below are allowed to exploit (some assumption for) the value of  $b$  when they take place to reduce the amount of work they require. As a result, from now on, I ignore the input space and its partition.

Physically, a ping-pong sequence is implemented by some continuous-time stochastic processes over  $W \times U$ . Any such process induces an associated discrete-time stochastic process over  $W \times U$ . That discrete-time process comprises a joint distribution  $Pr$  defined over a (possibly infinite) sequence of values  $(w_0, u_0), \dots, (w_t, u_t), (w_{t+1}, u_{t+1}), \dots$ . That distribution in turn induces a joint distribution over associated pairs of partition element labels,  $(w_0, u_0), \dots, (x_t, y_t), (x_{t+1}, y_{t+1}), \dots$ .

For calculational simplicity, I assume that  $\forall y \in Y$ , at the end of each stage in a ping-pong sequence that starts at any time  $t \in \mathbb{N}$ ,  $Pr(u | y)$  is the same distribution, which I write as  $q_{out}^y(u)$ . I make the analogous assumption for  $Pr(w | x)$  to define  $q_{proc}^x(w)$  (in addition to simplifying the analysis, this helps ensure that we are considering cyclic processes, a crucial issue whenever analyzing issues like the minimal amount of work needed to implement a desired map). Note that  $q_{out}^y(u) = 0$  if  $\mathcal{Y}(u) \neq y$ . To simplify the analysis further, I also assume that all “internal entropies”

of the processor macrostates are the same, *i.e.*,  $S(q_{out}^y(U))$  is independent of  $y$ , and similarly for the internal entropies of the output macrostates.

Also for calculational simplicity, I assume that at the end of each stage in a ping-pong sequence that starts at any time  $t \in \mathbb{N}$ , there is no interaction Hamiltonian coupling any of the three subsystems (though obviously, there must be such coupling at non-integer times). I also assume that at all such moments, the Hamiltonian over  $U$  is the same function, which I write as  $H_{out}$ . Therefore, for all such moments, the expected value of the Hamiltonian over  $U$  if the system is in state  $y_t$  at that time is:

$$\mathbb{E}(H_{out} | y) = \sum_u q_{out}^y(u) H_{out}(u) \tag{35}$$

Similarly,  $H_{in}$  and  $H_{proc}$  define the Hamiltonians at all such moments, over the input and processor subsystems, respectively.

I will refer to any quadruple  $(W, \mathcal{X}, U, \mathcal{Y})$  and three associated Hamiltonians as an *organism*.

For future use, note that for any iteration  $t \in \mathbb{N}$ , initial distribution  $\mathcal{P}'(x_1)$ , conditional distribution  $\pi(y | x)$  and halting subset  $Y_{halt} \subseteq Y$ ,

$$\begin{aligned} \mathcal{P}'(y_t \in Y_{halt}) &= \sum_{y_t} \mathcal{P}'(y_t) I(y_t \in Y_{halt}) \\ &= \sum_{x_t, y_t} \mathcal{P}'(x_t) \pi(y | x)|_{x=x_t, y=y_t} I(y_t \in Y_{halt}) \end{aligned} \tag{36}$$

$$\mathcal{P}'(y_t | y_t \in Y_{halt}) = \frac{\sum_{x_t} \mathcal{P}'(x_t) \pi(y | x)|_{x=x_t, y=y_t} I(y_t \in Y_{halt})}{\sum_{x_t, y_t} \mathcal{P}'(x_t) \pi(y | x)|_{x=x_t, y=y_t} I(y_t \in Y_{halt})} \tag{37}$$

and similarly:

$$\mathcal{P}'(x_{t+1} | y_t \notin Y_{halt}) = \frac{\sum_{x_t} \mathcal{P}'(x_t) \pi(y | x)|_{x=x_t, y=x_{t+1}} I(x_{t+1} \notin Y_{halt})}{\sum_{x_t, y_t} \mathcal{P}'(x_t) \pi(y | x)|_{x=x_t, y=x_{t+1}} I(x_{t+1} \notin Y_{halt})} \tag{38}$$

Furthermore,

$$S(\mathcal{P}_t(X)) = - \sum_x \mathcal{P}_t(x) \ln[\mathcal{P}_t(x)] \tag{39}$$

$$S(\mathcal{P}_{t+1}(X)) = - \sum_{x, y} \mathcal{P}_t(x) \pi(y | x) \ln \left[ \sum_{x'} \mathcal{P}_t(x') \pi(y | x') \right] \tag{40}$$

I end this subsection with some notational comments. I will sometimes abuse notation and put time indices on distributions rather than variables, *e.g.*, writing  $Pr_t(y)$  rather than  $Pr(y_t = y)$ . In addition, sometimes, I abuse notation with temporal subscripts. In particular, when the initial distribution over  $X$  is  $\mathcal{P}_1(x)$ , I sometimes use expressions like:

$$\mathcal{P}_t(w) \equiv \sum_x \mathcal{P}_t(x) q_{in}^x(w) \tag{41}$$

$$\mathcal{P}_t(u) \equiv \sum_y \mathcal{P}_t(y) q_{out}^y(u) \tag{42}$$

$$\mathcal{P}_t(y) \equiv \sum_{x_t} \mathcal{P}_t(x_t) \pi(y_t | x_t) \tag{43}$$

$$\mathcal{P}_{t+1}(x | y_t) \equiv \delta(x, y_t) \tag{44}$$

However, I will always be careful when writing joint distributions over variables from different moments of time, e.g., writing:

$$\begin{aligned} \mathcal{P}(y_{t+1}, x_t) &\equiv \mathcal{P}(y_{t+1} | x_t) \mathcal{P}(x_t) \\ &= \pi(y_{t+1} | x_t) \mathcal{P}_t(x_t) \end{aligned} \tag{45}$$

### 3.2. The Thermodynamics of Mapping an Input Space to an Output Space

Our goal is to construct a physical process  $\Lambda$  over an organism’s quadruple  $(W, \mathcal{X}, U, \mathcal{Y})$  that implements an iteration of a given ping-pong sequence above for any particular  $t$ . In addition, we want  $\Lambda$  to be thermodynamically optimal with the stipulated starting and ending joint Hamiltonians for all iterations of the ping-pong sequence when it is run on an initial joint distribution:

$$\mathcal{P}_1(x, y) = \mathcal{P}_1(x) \delta(y, 0) \tag{46}$$

In Appendix B, I present four separate GQ processes that implement stages (1), (2), (4) and (5) in a ping-pong sequence (and so implement the entire sequence). The GQ processes for stages (1), (4) and (5) are guaranteed to be thermodynamically reversible, for all  $t$ . However, each time- $t$  GQ process for stage (2) is parameterized by a distribution  $\mathcal{G}_t(x_t)$ . Intuitively, that distribution is a guess, made by the “designer” of the (time- $t$ ) stage (2) GQ process, for the marginal distribution over the values  $x_t$  at the beginning of the associated stage (1) GQ process. That stage (2) GQ process will also be thermodynamically reversible, if the distribution over  $x_t$  at the beginning of the stage (1) GQ process is in fact  $\mathcal{G}_t(x_t)$ . Therefore, for that input distribution, the sequence of GQ processes is thermodynamically optimal, as desired. However, as discussed below, in general, work will be dissipated if the stage (2) GQ process is applied when the distribution over  $x_t$  at the beginning of stage (1) differs from  $\mathcal{G}_t(x_t)$ .

I call such a sequence of five processes implementing an iteration of a ping-pong sequence an *organism process*. It is important to emphasize that I do *not* assume that any particular real biological system runs an organism process. An organism process provides a counterfactual model of how to implement a particular dynamics over  $X \times Y$ , a model that allows us to calculate the minimal work used by any actual biological system that implements that dynamics.

Suppose that an organism process always halts for any  $x_1$ , such that  $\mathcal{P}_1(x_1) \neq 0$ . Let  $\tau^*$  be the last iteration at which such an organism process may halt, for any of the inputs  $x_1$ , such that  $\mathcal{P}(x_1) \neq 0$  (note that if  $X$  is countably infinite,  $\tau^*$  might be countable infinity). Suppose further that no new input is received before  $\tau^*$  if the process halts at some  $\tau < \tau^*$  and that all microstates are constant from such a  $\tau$  up to  $\tau^*$  (so, no new work is done during such an interval). In light of the iterative nature of organism processes, this last assumption is equivalent to assuming that  $\pi(y_t | x_t) = \delta_{y_t, x_t}$  if  $x_t \in Y_{halt}$ .

I say that the organism process is *recursive* when all of these conditions are met, since that is the adjective used in the theory of Turing machines. For a recursive organism process, the ending distribution over  $y$  is:

$$\mathcal{P}(y_{\tau^*}) = \sum_{x_1, \dots, x_{\tau^*}} \pi(y_{\tau^*} | x_{\tau^*}) \mathcal{P}_1(x_1) \prod_{t=1}^{\tau^*} \pi(x_t | x_{t-1}) \tag{47}$$

and:

$$\mathcal{P}(y_{\tau^*} | x_1) = \sum_{x_2, \dots, x_{\tau^*}} \pi(y_{\tau^*} | x_{\tau^*}) \prod_{t=1}^{\tau^*} \pi(x_t | x_{t-1}) \tag{48}$$

**Proposition 2.** Fix any recursive organism process, iteration  $t \in \mathbb{N}$ , initial distributions  $\mathcal{P}_1(x), \mathcal{P}'_1(x)$ , conditional distribution  $\pi(y | x)$  and halting subset  $Y_{halt} \subseteq Y$ .

1. With probability  $\mathcal{P}'(y_t \in Y_{halt})$ , the ping-pong sequence at iteration  $t$  of the associated organism process maps the distribution:

$$\mathcal{P}'(x_t)\delta(y_{t-1}, 0) \rightarrow \delta(x_t, 0)\mathcal{P}'(y_t | y_t \in Y_{halt})$$

and then halts, and with probability  $1 - \mathcal{P}'(y_t \in Y_{halt})$ , it instead maps:

$$\mathcal{P}'(x_t)\delta(y_{t-1}, 0) \rightarrow \mathcal{P}(x_{t+1} | y_t \notin Y_{halt})\delta(y_t, 0)$$

and continues.

2. If  $\mathcal{G}_t = \mathcal{P}_t$  for all  $t \leq \tau^*$ , the total work the organism expends to map the initial distribution  $\mathcal{P}_1(x)$  to the ending distribution  $\mathcal{P}_{\tau^*}(y)$  is:

$$\begin{aligned} \Omega_{\mathcal{P}_1}^\pi &\equiv \sum_y \mathcal{P}_{\tau^*}(y)\mathbb{E}(H_{out} | y) - \mathbb{E}(H_{out} | y')|_{y'=0} - \sum_x \mathcal{P}_1(x)\mathbb{E}(H_{in} | x) + \mathbb{E}(H_{in} | x')|_{x'=0} \\ &\quad + kT(S(\mathcal{P}_1(X)) - S(\mathcal{P}_{\tau^*}(Y))) \end{aligned}$$

3. There is no physical process that both performs the same map as the organism process and that requires less work than the organism process does when applied to  $\mathcal{P}(x_t)\delta(y_t, 0)$ .

**Proof.** Repeated application of Proposition 1 gives the first result.

Next, combine Equation (71) in Appendix B, Equation (33) and our assumptions made just before Equation (35) to calculate the work needed to implement the GQ process of the first stage of an organism process at iteration  $t$ :

$$\begin{aligned} &\left[ \sum_{x,y,u} \left( \mathcal{P}_t(x)\pi(y | x)q_{out}^y(u) - q_{out}^0(u) \right) H_{out}(u) \right] - kT \left[ S(\mathcal{P}_t(Y)) - S(\mathcal{P}_{t-1}(Y)) \right] \\ &= \sum_y \mathcal{P}_t(y)\mathbb{E}(H_{out} | y) - \mathbb{E}(H_{out} | y')|_{y'=0} - kTS(\mathcal{P}_t(Y)) \end{aligned}$$

Analogous equations give the work for the remaining three GQ processes. Then, apply these equations repeatedly, starting with the distribution given in Equation (46) (note that all terms for iterations of the ping-pong sequence with  $t \in \{2, 3, \dots, \tau^* - 1\}$  cancel out). This gives the second result.

Finally, the third result is immediate from the assumption that  $\mathcal{G}_t = \mathcal{P}_t$  for all  $t$ , which guarantees that each iteration of the organism process is thermodynamically reversible.  $\square$

The first result in Proposition 2 means that no matter what the initial distribution over  $X$  is, the organism process updates that distribution according to  $\pi$ , halting whenever it produces a value in  $Y_{halt}$ . This is true even if the output of  $\pi$  depends on its input (as discussed in the Introduction, this property is violated for many of the physical processes considered in the literature).

The first terms in the definition of  $\Omega_{\mathcal{P}_1}^\pi$ , given by a sum of expected values of the Hamiltonian, can be interpreted as the “labor” done by the organism when processing  $x_1$  into  $y_{\tau^*}$ , e.g., by making and breaking chemical bonds. It quantifies the minimal amount of external free energy that must be used to implement the amount of labor that is (implicitly) specified by  $\pi$ . The remaining terms, a difference of entropies, represent the free energy required by the “computation” done by the organism when it undergoes  $\pi$ , independent of the labor done by the organism.

### 3.3. Input Distributions and Dissipated Work

Suppose that at the beginning of some iteration  $t$  of an organism process, the distribution over  $x_t$  is some  $\mathcal{P}(x_t)$  that differs from  $\mathcal{G}_t(x_t)$ , the prior distribution “built into” the (quenching Hamiltonians defining the) organism process. Then, as elaborated at the beginning of Section 3.2, in general, this iteration of the organism process will result in dissipated work.

As an example, such dissipation will occur if the organism process is used in an environment that generates inputs according to a distribution  $\mathcal{P}_1$  that differs from  $\mathcal{G}_0$ , the distribution “built into” the organism process. In the context of biology, if a biological system gets optimized by natural selection for one environment, but is then used in another one, it will necessarily operate (thermodynamically sub-optimally) in that second environment.

Note though that one could imagine designing an organism to operate optimally for a distribution over environments, since that is equivalent to a single average distribution over inputs. More precisely, a distribution  $Pr(\mathcal{P}_1)$  over environments is equivalent to a single environment generating inputs according to:

$$Pr(x_1) = \sum_{\mathcal{P}_1} Pr(\mathcal{P}_1)\mathcal{P}_1(x_1) \tag{49}$$

We can evaluate the thermodynamic cost  $\Omega_{\mathcal{P}_r}^\pi$  for this organism that behaves optimally for an uncertain environment.

As a comparison point, we can also evaluate the work used in an impossible scenario where  $\mathcal{P}_1$  varies stochastically, but the organism magically “knows” what each  $\mathcal{P}_1$  is before it receives an input sampled from that  $\mathcal{P}_1$ , and then changes its distributions  $\mathcal{G}_t$  accordingly to what the average thermodynamic cost in this impossible scenario would be

$$\sum_{\mathcal{P}_1} Pr(\mathcal{P}_1)\Omega_{\mathcal{P}_1}^\pi \tag{50}$$

In general

$$\Omega_{\mathcal{P}_r}^\pi \geq \sum_{\mathcal{P}_1} Pr(\mathcal{P}_1)\Omega_{\mathcal{P}_1}^\pi \tag{51}$$

with equality only if  $Pr(\cdot)$  is a delta function about one particular  $\mathcal{P}_1$ . So in general, even if an organism chooses its (fixed)  $\mathcal{G}_0$  to be optimal for an uncertain environment, it cannot do as well as it would if it could magically change  $\mathcal{G}_0$  appropriately before each new environment it encounters.

As a second example, in general, as one iterates an organism process, the initial distribution  $\mathcal{P}_1(x)$  is changed into a sequence of new distributions  $\{\mathcal{P}_1(x), \mathcal{P}_2(x), \dots\}$ . In general, many of these distributions will differ, *i.e.*, for many  $t'$ ,  $\mathcal{P}_{t'+1} \neq \mathcal{P}_{t'}$ . Accordingly, if one is using some particular physical device to implement the organism process, unless that device has a clock that it can use to update  $\mathcal{G}_t$  from one iteration to the next (to match the changes in  $\mathcal{P}_t$ ), the distribution  $\mathcal{G}_t$  built into the device will differ from  $\mathcal{P}_t$  at some times  $t$ . Therefore, without such a clock, work will be dissipated.

Bearing these caveats in mind, unless explicitly stated otherwise, in the sequel, I assume that the time- $t$  stage (2) GQ process of an organism makes the correct guess for the input distribution at the start of the time- $t$  ping-pong sequence, *i.e.*, that its parameter  $\mathcal{G}_t$  is always the same as the distribution over  $x$  at the beginning of the time- $t$  stage (1) process. In this case, the minimal free energy required by the organism is  $\Omega_{\mathcal{P}_1}^\pi$ , and no work is dissipated.

It is important to realize that in general, if one were to run a Q process over  $X$  in the second stage of an organism process, rather than a GQ process over  $X$  guided by  $Y$ , there would be nonzero dissipated work. The reason is that if we ran such a Q process, we would ignore the information in  $y_{t+1}$  concerning the variable we want to send to zero,  $x_t$ . In contrast, when we use a GQ process over  $X$  guided by  $Y$ , no information is ignored, and we maintain thermodynamic reversibility. The extra work of the Q process beyond that of the GQ process is:

$$kTS(X_t) - kTS(X_t | Y_{t+1}) = kTI(X_t; Y_{t+1}) \tag{52}$$

In other words, using the Q process would cause us to dissipate work  $kTI(X_t; Y_{t+1})$ . This amount of dissipated work equals zero if the output of  $\pi$  is independent of its input, as in bit erasure. It also

equals zero if  $P(x_t)$  is a delta function. However, for other  $\pi$  and  $P(x_t)$ , that dissipated work will be nonzero. In such situations, stage 2 would be thermodynamically irreversible if we used a Q process over  $X_t$  to set  $x$  to zero.

As a final comment, it is important to emphasize that no claim is being made that the only way to implement an organism process is with Q processes and/or GQ processes. However, the need to use the organism process in an appropriate environment, and for it to have a clock, should be generic, if we wish to avoid dissipated work.

### 3.4. Optimal Organisms

From now on, for simplicity, I restrict attention to recursive organism processes.

Recall that adding noise to  $\pi$  may reduce the amount of work required to implement it. Formally, Proposition 2 tells us that everything else being equal, the larger  $S(\mathcal{P}_{\tau^*}(Y))$  is, the less work is required to implement the associated  $\pi$  (indeed, the thermodynamically-optimal implementation of a one-to-many map  $\pi$  actually draws in free energy from the heat bath, rather than requiring free energy that ends up being dumped into that heat bath). This implies that an organism will want to implement a  $\pi$  that is as noisy as possible.

In addition, not all maps  $x_1 \rightarrow y_{\tau^*}$  are equally important to an organism’s reproductive fitness. It will be important to be very precise in what output is produced for some inputs  $x_1$ , but for other inputs, precision is not so important. Indeed, for some inputs, it may not matter at all what output the organism produces in response.

In light of this, natural selection would be expected to favor  $\pi$ ’s that are as noisy as possible, while still being precise for those inputs where reproductive fitness requires it. To simplify the situation, there are two contributions to the reproductive fitness of an organism that implements some particular  $\pi$ : the free energy (and other resources) required by that implementation and the “phenotypic fitness” that would arise by implementing  $\pi$  even if there were no resources required to implement it.

Therefore, there will be a tradeoff between the resource cost of being precise in  $\pi$  with the phenotypic fitness benefit of being precise. In particular, there will be a tradeoff between the thermodynamic cost of being precise in  $\pi$  (given by the minimal free energy that needs to be used to implement  $\pi$ ) and the phenotypic fitness of that  $\pi$ . In this subsection, I use an extremely simplified and abstracted model of reproductive fitness of an organism to determine what  $\pi$  optimizes this tradeoff.

To start, suppose we are given a real-valued *phenotypic fitness* function  $f(x_1, y_{\tau^*})$ . This quantifies the benefit to the organism of being precise in what output it produces in response to its inputs. More precisely,  $f(x_1, y_{\tau^*})$  quantifies the impact on the reproductive fitness of the organism that arises if it outputs  $y_{\tau^*}$  in response to an input  $x_1$  it received, minus the effect on reproductive fitness of how the organism generated that response. That second part of the definition means that behavioral fitness does not include energetic costs associated with mapping  $x_1 \rightarrow y_{\tau^*}$ . Therefore, it includes neither the work required to compute a map taking  $x_1 \rightarrow y_{\tau^*}$  nor the labor involved in carrying out that map going into  $f$  (note that in some toy models,  $f(x_1, y_{\tau^*})$  would be an expectation value of an appropriate quantity, taken over states of the environment, and conditioned on  $x_1$  and  $y_{\tau^*}$ ). For an input distribution  $\mathcal{P}_1(x)$  and conditional distribution  $\pi$ , expected phenotypic fitness is:

$$\mathbb{E}_{\mathcal{P}_1, \pi}(f) = \sum_{x_1, y_{\tau^*}} \mathcal{P}_1(x_1) \mathcal{P}(y_{\tau^*} | x_1) f(x_1, y_{\tau^*}) \tag{53}$$

where  $\mathcal{P}(y_{\tau^*} | x_1)$  is given by Equation (48).

The expected phenotypic fitness of an organism if it implements  $\pi$  on the initial distribution  $\mathcal{P}_1$  is only one contribution to the overall reproductive fitness of the organism. In addition, there is a reproductive fitness cost to the organism that depends on the specific physical process it uses to

implement  $\pi$  on  $\mathcal{P}_1$ . In particular, there is such a cost arising from the physical resources that the process requires.

There are several contributions to this cost. In particular, different physical processes for implementing  $\pi$  will require different sets of chemicals from the environment, will result in different chemical waste products, *etc.* Here, I ignore such “material” costs of the particular physical process the organism uses to implement  $\pi$  on  $\mathcal{P}_1$ .

However, in addition to these material costs of the process, there is also a cost arising from the thermodynamic work required to run that process. If we can use a thermodynamically-reversibly process, then by Equation (49), for fixed  $\mathcal{P}_1$  and  $\pi$ , the minimal possible such required work is  $\Omega_{\mathcal{P}_1}^\pi$ . Of course, in many biological scenarios, it is not possible to use a thermodynamically-reversible organism process to implement  $\pi$ . As discussed in Section 3.3, this is the case if the organism process is “designed” for an environment that generates inputs  $x$  according to  $\mathcal{G}_1(x)$  while the actual environment in which the process is used generates inputs according to some  $\mathcal{P}_1 \neq \mathcal{G}_1$ . However, there are other reasons why there might have to be non-zero dissipated work. In particular, there is non-zero dissipated work if  $\pi$  must be completed quickly, and so, it cannot be implemented using a quasi-static process (it does not do an impala any good to be able to compute the optimal direction in which to flee a tiger chasing it, if it takes the impala an infinite amount of time to complete that computation). Additionally, of course, it may be that a minimal amount of work must be dissipated simply because of the limited kinds of biochemical systems available to a real organism.

I make several different simplifying assumptions:

1. In some biological scenarios, the amount of such dissipated work that cannot be avoided in implementing  $\pi$ ,  $\hat{W}_{\mathcal{P}_1}^\pi$ , will be comparable to (or even dominate) the minimal amount of reversible work needed to implement  $\pi$ ,  $\Omega_{\mathcal{P}_1}^\pi$ . However, for simplicity, in the sequel, I concentrate solely on the dependence on  $\pi$  of the reproductive fitness of a process that implements  $\pi$  that arises due to its effect on  $W_{\mathcal{P}_1}^\pi$ . Equivalently, I assume that I can approximate differences  $\hat{W}_{\mathcal{P}_1}^\pi - \hat{W}_{\mathcal{P}_1}^{\pi'}$  as equal to  $\hat{W}_{\mathcal{P}_1}^\pi - \hat{W}_{\mathcal{P}_1}^{\pi'}$  up to an overall proportionality constant.
2. Real organisms have internal energy stores that allow them to use free energy extracted from the environment at a time  $t' < 1$  to drive a process at time  $t = 1$ , thereby “smoothing out” their free energy needs. For simplicity, I ignore such energy stores. Under this simplification, the organism needs to extract at least  $\Omega_{\mathcal{P}_1}^\pi$  of free energy from its environment to implement a single iteration of  $\pi$  on  $\mathcal{P}_1$ . That minimal amount of needed free energy is another contribution to the “reproductive fitness cost to the organism of physically implementing  $\pi$  starting from the input distribution  $\mathcal{P}_1$ ”.
3. As another simplifying assumption, I suppose that the (expected) reproductive fitness of an organism that implements the map  $\pi$  starting from  $\mathcal{P}_1$  is just:

$$\mathcal{F}(\mathcal{P}_1, \pi, f) \equiv \alpha \mathbb{E}_{\mathcal{P}_1, \pi}(f) - \Omega_{\mathcal{P}_1}^\pi \tag{54}$$

Therefore,  $\alpha$  is the benefit to the organism’s reproductive fitness of increasing  $f$  by one, measured in units of energy. This ignores all effects on the distribution  $\mathcal{P}_1$  that would arise by having different  $\pi$  implemented at times earlier than  $t = 1$ . It also ignores the possible impact on reproductive fitness of the organism’s implementing particular sequences of multiple  $y$ ’s (future work involves weakening all of these assumptions, with particular attention to this last one). Under this assumption, varying  $\pi$  has no effect on  $S(X_1)$ , the initial entropy over processor states. Similarly, it has no effect on the expected value of the Hamiltonian then.

Combining these assumptions with Proposition 2, we see that after removing all terms in  $\Omega_{\mathcal{P}_1}^\pi$  that do not depend on  $\pi$ , we are left with  $\sum_y \mathcal{P}_{\tau^*}(y) \mathbb{E}(H_{out} | y) - kTS(\mathcal{P}_{\tau^*}(Y))$ . This gives the following result:

**Corollary 1.** Given the assumptions discussed above, up to an additive constant that does not depend on  $\pi$ :

$$\mathcal{F}(\mathcal{P}_1, \pi, f) = \sum_{x_1, y_{\tau^*}} \mathcal{P}(x_1) \mathcal{P}(y_{\tau^*} | x_1) \left\{ \alpha f(x_1, y_{\tau^*}) - H_{out}(y_{\tau^*}) - kT \ln \left[ \sum_{x'_1} \mathcal{P}_1(x'_1) \mathcal{P}(y_{\tau^*} | x'_1) \right] \right\}$$

The first term in Corollary 1 reflects the impact of  $\pi$  on the phenotypic fitness of the organism. The second term reflects the impact of  $\pi$  on the amount of labor the organism does. Finally, the last term reflects the impact of  $\pi$  on the amount of computation the organism does; the greater the entropy of  $y_{\tau^*}$ , the less total computation is done. In different biological scenarios, the relative sizes of these three terms may change radically. In some senses, Corollary 1 can be viewed as an elaboration of [58], where the “cost of sensing” constant in that paper is decomposed into labor and computation costs.

From now on, for simplicity, I assume that  $Y_{halt} = Y$ . So no matter what the input is, the organism process runs  $\pi$  exactly once to produce the output. Returning to our actual optimization problem, by Lagrange multipliers, if the  $\pi$  that maximizes the expression in Corollary 1 lies in the interior of the feasible set, then it is the solution to a set of coupled nonlinear equations, one equation for each pair  $(x_1, y_1)$ :

$$\mathcal{P}(x_1) \left\{ H_{out}(y_1) - \alpha f(x_1, y_1) + kT \left( \ln \left[ \sum_{x'_1} \mathcal{P}(x'_1) \pi(y_1 | x'_1) \right] + 1 \right) \right\} = \lambda_{x_1} \tag{55}$$

where the  $\lambda_{x_1}$  are the Lagrange multipliers ensuring that  $\sum_{y_1} \pi(y_1 | x_1) = 1$  for all  $x_1 \in X$ . Unfortunately, in general, the solution may not lie in the interior, so that we have a non-trivial optimization problem.

However, suppose we replace the quantity:

$$- \sum_{x_1, y_1} \mathcal{P}_1(x_1) \pi(y_1 | x_1) \ln \left[ \sum_{x'_1} \mathcal{P}_1(x'_1) \pi(y_1 | x'_1) \right] = S(Y_1) \tag{56}$$

in Corollary 1 with  $S(Y_1 | X_1)$ . Since  $S(Y_1 | X_1) \leq S(Y_1)$  [50,51], this modification gives us a lower bound on expected reproductive fitness:

$$\hat{\mathcal{F}}(\mathcal{P}_1, \pi, f) \equiv \sum_{x_1, y_1} \mathcal{P}_1(x_1) \pi(y_1 | x_1) \left\{ \alpha f(x_1, y_1) - H_{out}(y_1) - kT \ln \left[ \pi(y_1 | x_1) \right] \right\} \tag{57}$$

$$\leq \mathcal{F}(\mathcal{P}_1, \pi, f) \tag{58}$$

The  $\pi$  that maximizes  $\hat{\mathcal{F}}(\mathcal{P}_1, \pi, f)$  is just a set of Boltzmann distributions:

$$\pi(y_1 | x_1) \propto \exp \left( \frac{\alpha f(x_1, y_1) - H_{out}(y_1)}{kT} \right) \tag{59}$$

For each  $x_1$ , this approximately optimal conditional distribution puts more weight on  $y_1$  if the associated phenotypic fitness is high, while putting less weight on  $y_1$  if the associated energy is large. In addition, we can use this distribution to construct a lower bound on the maximal value of the expected reproductive fitness:

**Corollary 2.** Given the assumptions discussed above,

$$\max_{\pi} \mathcal{F}(\mathcal{P}_1, \pi, f) \geq -kT \sum_{x_1} \mathcal{P}(x_1) \ln \left[ \sum_{y_1} \exp \left( \frac{\alpha f(x_1, y_1) - H_{out}(y_1)}{kT} \right) \right]$$

**Proof.** Write:

$$\begin{aligned} \hat{\mathcal{F}}(\mathcal{P}_1, \pi, f) &= \sum_{x_1, y_1} \mathcal{P}_1(x_1) \left( \pi(y_1 | x_1) \left\{ \alpha f(x_1, y_1) - H_{out}(y_1) - kT \ln \left[ \pi(y_1 | x_1) \right] \right\} \right) \\ &\equiv \sum_{x_1, y_1} \mathcal{P}_1(x_1) \hat{\mathcal{F}}(x_1, \pi, f) \end{aligned} \tag{60}$$

Each term  $\hat{\mathcal{F}}(x_1, \pi, f)$  in the summand depends on the  $Y$ -space distribution  $\pi(\cdot | x_1)$ , but no other terms in  $\pi$ . Therefore, we can evaluate each such term  $\hat{\mathcal{F}}(x_1, \pi, f)$  separately for its maximizing (Boltzmann) distribution  $\pi(\cdot | x_1)$ . In the usual way, this is given by the log of the associated partition function (normalization constant)  $z(x_1)$ , since for any  $x_1$  and associated Boltzmann  $\pi(\cdot | x_1)$ ,

$$\begin{aligned} S(Y_1 | x_1) &= - \sum_{y_1} \pi(y_1 | x_1) \ln[\pi(y_1 | x_1)] \\ &= - \sum_{y_1} \frac{\exp(\beta[\alpha f(x_1, y_1) - H_{out}(y_1)])}{z(x_1)} \ln \left[ \frac{\exp(\beta[\alpha f(x_1, y_1) - H_{out}(y_1)])}{z(x_1)} \right] \\ &= - \sum_{y_1} \pi(y_1 | x_1) (\beta[\alpha f(x_1, y_1) - H_{out}(y_1)]) - \ln[z(x_1)] \end{aligned} \tag{61}$$

where  $\beta \equiv 1/kT$ , as usual. Comparing to Equation (60) establishes that:

$$\hat{\mathcal{F}}(x_1, \pi, f) = -kT \ln[z(x_1)] \tag{62}$$

and then gives the claimed result.  $\square$

As an aside, suppose we had  $X = Y$ ,  $f(x, x) = 0$  for all  $x$  and that  $f$  were non-negative. Then if in addition the amount of expected work were given by the mutual information between  $X_1$  and  $Y_1$  rather than the difference in their entropies, our optimization problem would reduce to finding a point on the rate-distortion curve of conventional information theory, with  $f$  being the distortion function [51]. (See also [5] for a slight variant of rate-distortion theory, appropriate when  $Y$  differs from  $X$ , and so the requirement that  $f(x, x) = 0$  is dropped.) However as shown above the expected work to implement  $\pi$  does not depend on the precise coupling between  $x_1$  and  $y_1$  under  $\pi$ , but only the associated marginal distributions. So rate-distortion theory does not directly apply.

On the other hand, some of the same kinds of analysis used in rate-distortion theory can also be applied here. In particular, for any particular component  $\pi(y_1 | x_1)$  where  $\mathcal{P}_1(x_1) \neq 0$ , since  $\tau^* = 1$ ,

$$\begin{aligned} \frac{\partial^2}{\partial \pi(y_1 | x_1)^2} \mathcal{F}(x_1, \pi, f) &= \frac{\mathcal{P}_1(x_1)}{\mathcal{P}_1(y_1)} \\ &> 0 \end{aligned} \tag{63}$$

(where  $\mathcal{P}(y_1) = \sum_{x'_1} \mathcal{P}(x'_1) \pi(y_1 | x'_1)$ , as usual). So  $\mathcal{F}(x_1, \pi, f)$  is concave in every component of  $\pi$ . This means that the optimizing channel  $\pi$  may lie on the edge of the feasible region of conditional distributions. Note though that even if the solution is on the edge of the feasible region, in general for different  $x_1$  that optimize  $\pi(y_1 | x_1)$  will put all its probability mass on different edges of the unit simplex over  $Y$ . So when those edges are averaged under  $\mathcal{P}_1(x_1)$ , the result is a marginal distribution  $\mathcal{P}(y_1)$  that lies in the interior of the unit simplex over  $Y$ .

As a cautionary note, often in the real world, there is an inviolable upper bound on the rate at which a system can “harvest” free energy from its environment, *i.e.*, on how much free energy it can harvest per iteration of  $\pi$  (for example, a plant with a given surface area cannot harvest free energy at a faster rate than sunlight falls upon its surface). In that case, we are not interested in optimizing a quantity like  $\mathcal{F}(\mathcal{P}_1, \pi, f)$ , which is a weighted average of minimal free energy and expected phenotypic fitness per iteration of  $\pi$ . Instead, we have a constrained optimization problem with an

inequality constraint: find the  $\pi$  that maximizes some quantity (e.g., expected phenotypic fitness), subject to an inequality constraint on the free energy required to implement that  $\pi$ . Calculating solutions to these kinds of constrained optimization problem is the subject of future work.

#### 4. General Implications for Biology

Any work expended on an organism must first be acquired as free energy from the organism's environment. However, in many situations, there is a limit on the flux of free energy through an organism's immediate environment. Combined with the analysis above, such limits provide upper bounds on the "rate of (potentially noisy) computation" that can be achieved by a biological organism in that environment, once all energetic costs for the organism's labor (*i.e.*, its moving, making/breaking chemical bonds, *etc.*) are accounted for.

As an example, human brains do little labor. Therefore, these results bound the rate of computation of a human brain. Given the fitness cost of such computation (the brain uses  $\sim 20\%$  of the calories used by the human body), this bound contributes to the natural selective pressures on humans (in the limit that operational inefficiencies of the brain have already been minimized). In other words, these bounds suggest that natural selection imposes a tradeoff between the fitness quality of a brain's decisions and how much computation is required to make those decisions. In this regard, it is interesting to note that the brain is famously noisy, and as discussed above, noise in computation may reduce the total thermodynamic work required (see [6,10,59] for more about the energetic costs of the human brain and its relation to Landauer's bound).

As a second example, the rate of solar free energy incident upon the Earth provides an upper bound on the rate of computation that can be achieved by the biosphere (this bound holds for any choice for the partition of the biosphere's fine-grained space into macrostates, such that the dynamics over those macrostates executes  $\pi$ ). In particular, it provides an upper bound on the rate of computation that can be achieved by human civilization, if we remain on the surface of the Earth and only use sunlight to power our computation.

Despite the use of the term "organism", the analysis above is not limited to biological individuals. For example, one could take the input to be a current generation population of individuals, together with attributes of the environment shared by those individuals. We could also take the output to be the next generation of that population, after selective winnowing based on the attributes of the environment (e.g., via replicator dynamics). In this example, the bounds above do not refer to the "computation" performed by an individual, but rather by an entire population subject to natural selection. Therefore, those bounds give the minimal free energy required to run natural selection.

As a final example, one can use these results to analyze how the thermodynamic behavior of the biosphere changes with time. In particular, if one iterates  $\pi$  from one  $t$  to the next, then the associated initial distributions  $\mathcal{P}_t$  change. Accordingly, the minimal amount of free energy required to implement  $\pi$  changes. In theory, this allows us to calculate whether the rate of free energy required by the information processing of the terrestrial biosphere increases with time. Prosaically, has the rate of computation of the biosphere increased over evolutionary timescales? If it has done so for most of the time that the biosphere has existed, then one could plausibly view the fraction of free energy flux from the Sun that the biosphere uses as a measure of the "complexity" of the biosphere, a measure that has been increasing throughout the lifetime of the biosphere.

Note as well that there is a fixed current value of the total free energy flux incident on the biosphere (from both sunlight and, to a much smaller degree, geologic processes). By the results presented above, this rate of free energy flux gives an upper bound on the rate of computation that humanity as a whole can ever achieve, if it monopolizes all resources of Earth, but restricts itself to the surface of Earth.

## 5. Discussion

The noisier the input-output map  $\pi$  of a biological organism, the less free energy the organism needs to acquire from its environment to implement that map. Indeed, by using a sufficiently noisy  $\pi$ , an organism can *increase* its stored free energy. Therefore, noise might not just be a hindrance that an organism needs to circumvent; an organism may actually exploit noise, to “recharge its battery”.

In addition, not all maps  $x_t \rightarrow y_{t+1}$  are equally important to an organism’s reproductive fitness. In light of this, natural selection would be expected to favor  $\pi$ ’s that are as noisy as possible, while still being precise for those inputs where reproductive fitness requires it.

In this paper, I calculated what  $\pi$  optimizes this tradeoff. This calculation provides insight into what phenotypes natural selection might be expected to favor. Note though that in the real world, there are many other thermodynamic factors that are important in addition to the cost of processing sensor readings (inputs) into outputs (actions). For example, there are the costs of acquiring the sensor information in the first place and of internal storage of such information, for future use. Moreover, in the real world, sensor readings do not arrive in an i.i.d. basis, as assumed in this paper. Indeed, in real biological systems, often, the current sensor reading, reflecting the recent state of the environment, reflects previous actions by the organism that affected that same environment (in other words, real biological organisms often behave like feedback controllers). All of these effects would modify the calculations done in this paper.

In addition, in the real world, there are strong limits on how much time a biological system can take to perform its computations, physical labor and rearranging of matter, due to environmental exigencies (simply put, if the biological system is not fast enough, it may be killed). These temporal constraints mean that biological systems cannot use fully reversible thermodynamics. Therefore, these temporal constraints increase the free energy required for the biological system to perform computation, labor and/or rearrangement of matter.

Future work involves extending the analysis of this paper to account for such thermodynamic effects. Combined with other non-thermodynamic resource restrictions that real biological organisms face, such future analysis should help us understand how closely the organisms that natural selection has produced match the best ones possible.

**Acknowledgment:** I would like to thank Daniel Polani, Sankaran Ramakrishnan and especially Artemy Kolchinsky for many helpful discussions and the Santa Fe Institute for helping to support this research. This paper was made possible through the support of Grant No. TWCF0079/AB47 from the Templeton World Charity Foundation and Grant No. FQXi-RH13-1349 from the FQXi foundation. The opinions expressed in this paper are those of the author and do not necessarily reflect the view of Templeton World Charity Foundation.

**Conflicts of Interest:** The author declares no conflict of interest. The funding sponsors had no role in the design of the study; in the collection, analyses or interpretation of data; in the writing of the manuscript; nor in the decision to publish the results.

## Appendix A: Proof of Proposition 1

We begin with the following lemma:

**Lemma 1.** A GQ process over  $R$  guided by  $V$  (for conditional distribution  $\pi$  and initial distribution  $\rho^t(r, s)$ ) will transform any initial distribution:

$$p^t(r, s) = \sum_v p^t(v) \rho^t(s | v) p^t(r | v) \quad (64)$$

into a distribution:

$$p^{t+1}(r, s) = \sum_v p^t(v) \rho^t(s | v) \pi(r | v) \quad (65)$$

**Proof.** Fix some  $v^*$  by sampling  $p^t(v)$ . Since in a GQ, microstates only change during the quasi-static relaxation, after the first quench,  $s$  and, therefore,  $v$  still equal  $v^*$ . Due to the infinite potential barriers in  $\mathcal{S}$ , while  $s$  may change during that relaxation,  $v$  will not, and so,  $v^{t+1} = v^* = v_t$ . Therefore:

$$H_{quench;int}^t(r, s) \equiv -kT \ln[\pi(r | v_t)] \tag{66}$$

Now, at the end of the relaxation step,  $\rho(r, s)$  has settled to thermal equilibrium within the region  $R \times v_t \subset R \times V$ . Therefore, combining Equation (66) with Equations (29) and (28), we see that the distribution at the end of the relaxation is:

$$\begin{aligned} \rho^{t+1}(r, s) &\propto \exp\left(\frac{-H_{quench}^{t+1}(r, s)}{kT}\right) \delta(V(s), v_t) \\ &= \exp(\ln[\pi(r | v_t)] + \ln[\rho^t(s)]) \delta(V(s), v_t) \\ &= \pi(r | v_t) \rho^t(s) \delta(V(s), v_t) \\ &\propto \pi(r | v_t) \rho^t(s | v) \end{aligned} \tag{67}$$

Normalizing,

$$\rho^{t+1}(r, s) = \pi(r | v_t) \rho^t(s | v) \tag{68}$$

Averaging over  $v_t$  then gives  $p^{t+1}(r, s)$ :

$$p^{t+1}(r, s) = \sum_v p^t(v) \rho^t(s | v) \pi(r | v) \tag{69}$$

□

Next, note that  $\rho^t(s | v) = 0$  if  $s \notin V(s)$ . Therefore, if Equation (65) holds and we sum  $p^{t+1}(r, s)$  over all  $s \in V^{-1}(v)$  for an arbitrary  $v$ , we get:

$$p^{t+1}(r, v) = p^t(v) \pi(r | v) \tag{70}$$

Furthermore, no matter what  $\rho^t(s | v)$  is,  $p^t(r, v) = p^t(v) p^t(r | v)$ . As a result, Lemma 1 implies that a GQ process over  $R$  guided by  $V$  (for conditional distribution  $\pi$  and initial distribution  $\rho^t(r, s)$ ) will transform any initial distribution  $p^t(v) p^t(r | v)$  into a distribution  $p^t(v) \pi(r | v)$ . This is true whether or not  $p^t(v) = \rho^t(v)$  or  $p^t(r | v) = \rho^t(r | v)$ . This establishes the claim of Proposition 1 that the first “crucial feature” of GQ processes holds.

### Appendix B: The GQ Processes Iterating a Ping-Pong Sequence

In this section, I present the separate GQ processes for implementing the stages of a ping-pong sequence.

First, recall our assumption from just below the definition of a ping-pong sequence that at the end of any of its stages,  $Pr(u | y)$  is always the same distribution  $q_{out}^y(u)$  (and similarly for distributions like  $Pr(w | x)$ ). Accordingly, at the end of any stage of a ping-pong sequence that implements a GQ process over  $U$  guided by  $X$ , we can uniquely recover the conditional distribution  $Pr(u | x)$  from  $Pr(y | x)$ :

$$\bar{\pi}(u | x) \equiv \sum_y \pi(y | x) q_{out}^y(u) \tag{71}$$

(and similarly, for a GQ process over  $W$  guided by  $Y$ ). Conversely, we can always recover  $Pr(y | x)$  from  $Pr(u | x)$ , simply by marginalizing. Therefore, we can treat any distribution  $\bar{\pi}(u | x)$  defining

such a GQ process interchangeably with a distribution  $\pi(y | x)$  (and similarly, for distributions  $\bar{\pi}(w | y)$  and  $\pi(x | y)$  occurring in GQ processes over  $W$  guided by  $Y$ ).

1. To construct the GQ process for the first stage, begin by writing:

$$\begin{aligned} \rho^t(w, u) &= \sum_{x,y} \mathcal{G}_t(x) \delta(y, 0) q_{proc}^x(w) q_{out}^y(u) \\ &= q_{out}^0(u) \mathcal{G}_t(\mathcal{X}(w)) q_{proc}^{\mathcal{X}(w)}(w) \end{aligned} \tag{72}$$

where  $\mathcal{G}_t(x)$  is an assumption for the initial distribution over  $x$ , one that in general may be wrong. Furthermore, define the associated distribution:

$$\begin{aligned} \rho^t(u | x) &= \frac{\sum_{w \in \mathcal{X}(x)} \rho^t(w, u)}{\sum_{w' \in \mathcal{X}(x)} \rho^t(w, u')} \\ &= q_{out}^0(u) \end{aligned} \tag{73}$$

By Corollary 1, running a GQ process over  $Y$  guided by  $X$  for conditional distribution  $\bar{\pi}(u | x_t)$  and initial distribution  $\rho^t(w, u)$  will send any initial distribution  $\mathcal{P}_t(x) \rho^t(u | x) = \mathcal{P}_t(x) q_{out}^0(u)$  to a distribution  $\mathcal{P}_t(x) \bar{\pi}(u | x)$ . Therefore, in particular, it will send any initial  $x \rightarrow \bar{\pi}(u | x)$ . Due to the definition of  $q_{out}^y$  and Equation (71), the associated conditional distribution over  $y$  given  $x$ ,  $\sum_{u \in \mathcal{Y}(y)} \bar{\pi}(u | x)$ , is equal to  $\pi(y | x)$ . Accordingly, this GQ process implements the first stage of the organism process, as desired. In addition, it preserves the validity of our assumptions that  $Pr(u | y) = q_{out}^y(u)$  and similarly for  $Pr(w | x)$ .

Next, by the discussion at the end of Section 2.4, this GQ process will be thermodynamically reversible since by assumption,  $\rho^t(u | x)$  is the actual initial distribution over  $u$  conditioned on  $x$ .

2. To construct the GQ process for the second stage, start by defining an initial distribution based on a (possibly counterfactual) prior  $\mathcal{G}_t(x)$ :

$$\hat{\rho}(w_t, u_t) \equiv \sum_{x,y} \mathcal{G}_t(x) q_{proc}^x(w_t) \pi(y | x) q_{out}^y(u_t) \tag{74}$$

and the associated conditional distribution:

$$\hat{\rho}(w_t | y_t) = \frac{\sum_{u_t \in \mathcal{Y}(y_t)} \hat{\rho}(w_t, u_t)}{\sum_{w', u' \in \mathcal{Y}(y_t)} \hat{\rho}(w', u')} \tag{75}$$

Note that:

$$\hat{\rho}(w_t | y_t) = \mathcal{G}_t(x_t | y_t) q_{proc}^{x_t}(w_t) \tag{76}$$

where:

$$\mathcal{G}_t(x_t | y_t) \equiv \frac{\pi(y_t | x_t) \mathcal{G}_t(x_t)}{\sum_{x'} \pi(y_t | x') \mathcal{G}_t(x')} \tag{77}$$

Furthermore, define a conditional distribution:

$$\bar{\pi}(w_t | y_t) \equiv I(w_t \in \mathcal{X}(0)) q_{proc}^0(w_t) \tag{78}$$

Consider a GQ process over  $W$  guided by  $Y$  for conditional distribution  $\bar{\pi}(w_t | y_t)$  and initial distribution  $\hat{\rho}(w_t, u_t)$ . By Corollary 1, this GQ process implements the second stage, as desired. In addition, it preserves the validity of our assumptions that  $Pr(u | y) = q_{out}^y(u)$  and similarly fo  $Pr(w | x)$ .

- Next, by the discussion at the end of Section 2.4, this GQ process will be thermodynamically reversible if  $\hat{\rho}(w_t | y_{t+1})$  is the actual distribution over  $w_t$  conditioned on  $y_{t+1}$ . By Equation (77), this in general requires that  $\mathcal{G}_t(x_t)$ , the assumption for the initial distribution over  $x_t$  that is built into the step (ii) GQ process, is the actual initial distribution over  $x_t$ . As discussed at the end of Section 2.3, work will be dissipated if this is not the case. Physically, this means that if the device implementing this GQ process is thermodynamically optimal for one input distribution, but used with another, then work will be dissipated (the amount of work dissipated is given by the change in the Kullback–Leibler divergence between  $G$  and  $\mathcal{P}$  in that stage (4) GQ process; see [46]).
3. We can also implement the fourth stage by running a (different) GQ process over  $X$  guided by  $Y$ . This GQ process is a simple copy operation, *i.e.*, implements a single-valued, invertible function from  $y_{t+1}$  to the initialized state  $x$ . Therefore, it is thermodynamically reversible. Finally, we can implement the fifth stage by running an appropriate GQ process over  $Y$  guided by  $X$ . This process will also be thermodynamically reversible.

## References

1. Frank, S.A. Natural selection maximizes Fisher information. *J. Evolut. Biol.* **2009**, *22*, 231–244.
2. Frank, S.A. Natural selection. V. How to read the fundamental equations of evolutionary change in terms of information theory. *J. Evolut. Biol.* **2012**, *25*, 2377–2396.
3. Donaldson-Matasci, M.C.; Bergstrom, C.T.; Lachmann, M. The fitness value of information. *Oikos* **2010**, *119*, 219–230.
4. Krakauer, D.C. Darwinian demons, evolutionary complexity, and information maximization. *Chaos Interdiscip. J. Nonlinear Sci.* **2011**, *21*, 037110.
5. Taylor, S.F.; Tishby, N.; Bialek, W. Information and fitness. **2007**, arXiv:0712.4382.
6. Bullmore, E.; Sporns, O. The economy of brain network organization. *Nat. Rev. Neurosci.* **2012**, *13*, 336–349.
7. Sartori, P.; Granger, L.; Lee, C.F.; Horowitz, J.M. Thermodynamic costs of information processing in sensory adaptation. *PLoS Comput. Biol.* **2014**, *10*, e1003974.
8. Mehta, P.; Schwab, D.J. Energetic costs of cellular computation. *Proc. Natl. Acad. Sci. USA* **2012**, *109*, 17978–17982.
9. Mehta, P.; Lang, A.H.; Schwab, D.J. Landauer in the age of synthetic biology: Energy consumption and information processing in biochemical networks. *J. Stat. Phys.* **2015**, *162*, 1153–1166.
10. Laughlin, S.B. Energy as a constraint on the coding and processing of sensory information. *Curr. Opin. Neurobiol.* **2001**, *11*, 475–480.
11. Govern, C.C.; ten Wolde, P.R. Energy dissipation and noise correlations in biochemical sensing. *Phys. Rev. Lett.* **2014**, *113*, 258102.
12. Govern, C.C.; ten Wolde, P.R. Optimal resource allocation in cellular sensing systems. *Proc. Natl. Acad. Sci. USA* **2014**, *111*, 17486–17491.
13. Lestas, I.; Vinnicombe, G.; Paulsson, J. Fundamental limits on the suppression of molecular fluctuations. *Nature* **2010**, *467*, 174–178.
14. England, J.L. Statistical physics of self-replication. *J. Chem. Phys.* **2013**, *139*, 121923.
15. Landenmark, H.K.; Forgan, D.H.; Cockell, C.S. An estimate of the total DNA in the biosphere. *PLoS Biol.* **2015**, *13*, e1002168.
16. Landauer, R. Irreversibility and heat generation in the computing process. *IBM J. Res. Dev.* **1961**, *5*, 183–191.
17. Landauer, R. Minimal energy requirements in communication. *Science* **1996**, *272*, 1914–1918.
18. Landauer, R. The physical nature of information. *Physics Lett. A* **1996**, *217*, 188–193.
19. Bennett, C.H. Logical reversibility of computation. *IBM J. Res. Dev.* **1973**, *17*, 525–532.
20. Bennett, C.H. The thermodynamics of computation—A review. *Int. J. Theor. Phys.* **1982**, *21*, 905–940.
21. Bennett, C.H. Time/space trade-offs for reversible computation. *SIAM J. Comput.* **1989**, *18*, 766–776.
22. Bennett, C.H. Notes on Landauer’s principle, reversible computation, and Maxwell’s Demon. *Stud. Hist. Philos. Sci. B* **2003**, *34*, 501–510.
23. Maroney, O. Generalizing Landauer’s principle. *Phys. Rev. E* **2009**, *79*, 031105.

24. Plenio, M.B.; Vitelli, V. The physics of forgetting: Landauer's erasure principle and information theory. *Contemp. Phys.* **2001**, *42*, 25–60.
25. Shizume, K. Heat generation required by information erasure. *Phys. Rev. E* **1995**, *52*, 3495–3499.
26. Fredkin, E.; Toffoli, T. *Conservative Logic*; Springer: Berlin/Heidelberg, Germany, 2002.
27. Faist, P.; Dupuis, F.; Oppenheim, J.; Renner, R. A quantitative Landauer's principle. **2012**, arXiv:1211.1037.
28. Touchette, H.; Lloyd, S. Information-theoretic approach to the study of control systems. *Physica A* **2004**, *331*, 140–172.
29. Sagawa, T.; Ueda, M. Minimal energy cost for thermodynamic information processing: Measurement and information erasure. *Phys. Rev. Lett.* **2009**, *102*, 250602.
30. Dillenschneider, R.; Lutz, E. Comment on "Minimal Energy Cost for Thermodynamic Information Processing: Measurement and Information Erasure". *Phys. Rev. Lett.* **2010**, *104*, 198903.
31. Sagawa, T.; Ueda, M. Fluctuation theorem with information exchange: Role of correlations in stochastic thermodynamics. *Phys. Rev. Lett.* **2012**, *109*, 180602.
32. Crooks, G.E. Entropy production fluctuation theorem and the nonequilibrium work relation for free energy differences. *Phys. Rev. E* **1999**, *60*, 2721–2726.
33. Crooks, G.E. Nonequilibrium measurements of free energy differences for microscopically reversible Markovian systems. *J. Stat. Phys.* **1998**, *90*, 1481–1487.
34. Janna, F.C.; Moukalled, F.; Gómez, C.A. A Simple Derivation of Crooks Relation. *Int. J. Thermodyn.* **2013**, *16*, 97–101.
35. Jarzynski, C. Nonequilibrium equality for free energy differences. *Phys. Rev. Lett.* **1997**, *78*, doi:10.1103/PhysRevLett.78.2690.
36. Esposito, M.; van den Broeck, C. Second law and Landauer principle far from equilibrium. *Europhys. Lett.* **2011**, *95*, 40004.
37. Esposito, M.; van den Broeck, C. Three faces of the second law. I. Master equation formulation. *Phys. Rev. E* **2010**, *82*, 011143.
38. Parrondo, J.M.; Horowitz, J.M.; Sagawa, T. Thermodynamics of information. *Nat. Phys.* **2015**, *11*, 131–139.
39. Pollard, B.S. A Second Law for Open Markov Processes. **2014**, arXiv:1410.6531.
40. Seifert, U. Stochastic thermodynamics, fluctuation theorems and molecular machines. *Rep. Prog. Phys.* **2012**, *75*, 126001.
41. Takara, K.; Hasegawa, H.H.; Driebe, D. Generalization of the second law for a transition between nonequilibrium states. *Phys. Lett. A* **2010**, *375*, 88–92.
42. Hasegawa, H.H.; Ishikawa, J.; Takara, K.; Driebe, D. Generalization of the second law for a nonequilibrium initial state. *Phys. Lett. A* **2010**, *374*, 1001–1004.
43. Prokopenko, M.; Einav, I. Information thermodynamics of near-equilibrium computation. *Phys. Rev. E* **2015**, *91*, 062143.
44. Sagawa, T. Thermodynamic and logical reversibilities revisited. *J. Stat. Mech.* **2014**, *2014*, P03025.
45. Mandal, D.; Jarzynski, C. Work and information processing in a solvable model of Maxwell's demon. *Proc. Natl. Acad. Sci. USA* **2012**, *109*, 11641–11645.
46. Wolpert, D.H. Extending Landauer's bound from bit erasure to arbitrary computation. **2015**, arXiv:1508.05319.
47. Barato, A.C.; Seifert, U. Stochastic thermodynamics with information reservoirs. *Phys. Rev. E* **2014**, *90*, 042150.
48. Deffner, S.; Jarzynski, C. Information processing and the second law of thermodynamics: An inclusive, Hamiltonian approach. *Phys. Rev. X* **2013**, *3*, 041003.
49. Barato, A.C.; Seifert, U. An autonomous and reversible Maxwell's demon. *Europhys. Lett.* **2013**, *101*, 60001.
50. Mackay, D. *Information Theory, Inference, and Learning Algorithms*; Cambridge University Press: Cambridge, UK, 2003.
51. Cover, T.; Thomas, J. *Elements of Information Theory*; Wiley: New York, NY, USA, 1991.
52. Yeung, R.W. *A First Course in Information Theory*; Springer: Berlin/Heidelberg, Germany, 2012.
53. Reif, F. *Fundamentals of Statistical and Thermal Physics*; McGraw-Hill: New York, NY, USA, 1965.
54. Still, S.; Sivak, D.A.; Bell, A.J.; Crooks, G.E. Thermodynamics of prediction. *Phys. Rev. Lett.* **2012**, *109*, 120604.
55. Hopcroft, J.E.; Motwani, R.; Ullman, J.D. *Introduction to Automata Theory, Languages and Computability*; Addison-Wesley: Boston, MA, USA, 2000.

56. Li, M.; Vitányi, P. *An Introduction to Kolmogorov Complexity and Its Applications*; Springer: Berlin/Heidelberg, Germany, 2008.
57. Grunwald, P.; Vitányi, P. Shannon information and Kolmogorov complexity. **2004**, arXiv:cs/0410002.
58. Kussell, E.; Leibler, S. Phenotypic diversity, population growth, and information in fluctuating environments. *Science* **2005**, *309*, 2075–2078.
59. Sandberg, A. Energetics of the brain and AI. **2016**, arXiv:1602.04019.



© 2016 by the author; licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons by Attribution (CC-BY) license (<http://creativecommons.org/licenses/by/4.0/>).