# New Journal of Physics

The open access journal at the forefront of physics

CrossMark

**PAPER**

# Total cost of operating an information engine

Jaegon Um[1,2], Haye Hinrichsen[1,5], Chulan Kwon[3] and Hyunggyu Park[4]

[1] Universität Würzburg, Fakultät für Physik und Astronomie, 97074 Würzburg, Germany
[2] Quantum Universe Center, Korea Institute for Advanced Study, Seoul 130-722, Korea
[3] Department of Physics, Myongji University, Yongin 449-728, Korea
[4] School of Physics, Korea Institute for Advanced Study, Seoul 130-722, Korea
[5] Author to whom any correspondence should be addressed.

E-mail: hinrichsen@physik.uni-wuerzburg.de

## Abstract

We study a two-level system controlled in a discrete feedback loop, modeling both the system and the controller in terms of stochastic Markov processes. We find that the extracted work, which is known to be bounded from above by the mutual information acquired during measurement, has to be compensated by an additional energy supply during the measurement process itself, which is bounded by the same mutual information from below. Our results confirm that the total cost of operating an information engine is in full agreement with the conventional second law of thermodynamics. We also consider the efficiency of the information engine as a function of the cycle time and discuss the operating condition for maximal power generation. Moreover, we find that the entropy production of our information engine is maximal for maximal efficiency, in sharp contrast to conventional reversible heat engines.
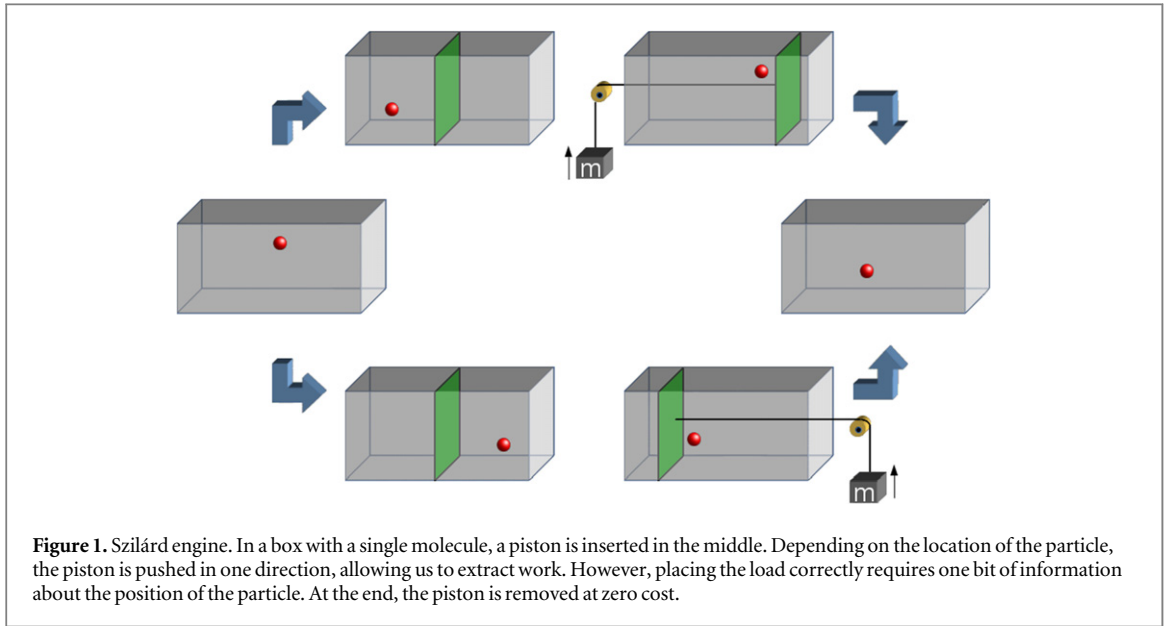
## 1. Introduction

In 1867, Maxwell created a thought experiment to demonstrate a possible violation of the second law of thermodynamics: a thermally isolated container with a gas is divided into two parts, and a fictitious demon opens or closes a door between the two parts depending on the velocity of the approaching particles, creating an increase of the temperature in one of the compartments [1].

Smoluchowski was the first to provide an explanation why Maxwell's demon does not work. To this end, he modeled the demon mechanically by a trapdoor combined with a gentle spring. The trapdoor acted as a valve in such a way that fast particles coming from one side can open the door while slow ones cannot, leading to a pressure difference. However, taking the full dynamics of the apparatus into account, Smoluchowski demonstrated that the energy of the spring system itself equilibrates at such a high energy that it opens and closes essentially randomly, leading to the same pressure difference as if the trapdoor was always open.

In 1928, Szilárd refined the concept of Maxwell's demon, suggesting what is known today as the *Szilárd engine* [3]. The starting point is a box that contains only one particle (see figure 1). If a wall is inserted in the middle, the particle will be in one of the two parts. Expanding the volume isothermally to its original size by moving the shutter into the empty half of the box, one can extract the work $W = k_B T \ln 2$. However, this requires knowing in which compartment the particle actually is, demonstrating that the possession of information can be converted into physical work. Thus, in order to keep a Szilárd engine running, a closed loop of measurement and feedback is needed. Very recently, such a feedback scheme could be realized experimentally for the first time [4].

Since the feedback mechanism in the Szilárd engine was perceived as information processing at zero cost, it seemed to produce usable work from nothing, violating the second law. However, as shown by Landauer [5] in 1961, any irreversible logical operation requires one to apply a well-defined minimum of work. In the case of the Szilárd engine, the demon has to store the information about the particle position in a single bit. According to Landauer's principle, resetting this bit requires one to convert a work of $k_B T \ln 2$ into heat. As shown by Bennett

**Figure 1.** Szilárd engine. In a box with a single molecule, a piston is inserted in the middle. Depending on the location of the particle, the piston is pushed in one direction, allowing us to extract work. However, placing the load correctly requires one bit of information about the position of the particle. At the end, the piston is removed at zero cost.

[6], who realized measurement and feedback as reversible processes, this extra work for resetting restores the second law.

Recently Maxwell's demon attracted renewed attention, as it was shown that a system in a feedback loop obeys an integral fluctuation theorem (IFT) of the form

$$\langle e^{\beta W_{ex} - \Delta I} \rangle = 1, \tag{1}$$

implying $\beta \langle W_{ex} \rangle \leqslant \langle \Delta I \rangle$, where $W_{ex}$ is the extracted work during feedback, $\beta$ is the inverse temperature, and $\Delta I$ is the gained mutual information between system and demon during the measurement [7–9]. This fluctuation theorem implies that in each cycle of the engine the extracted work is limited by the gained mutual information —a highly plausible result that nicely demonstrates the equivalence of thermodynamic work and information.
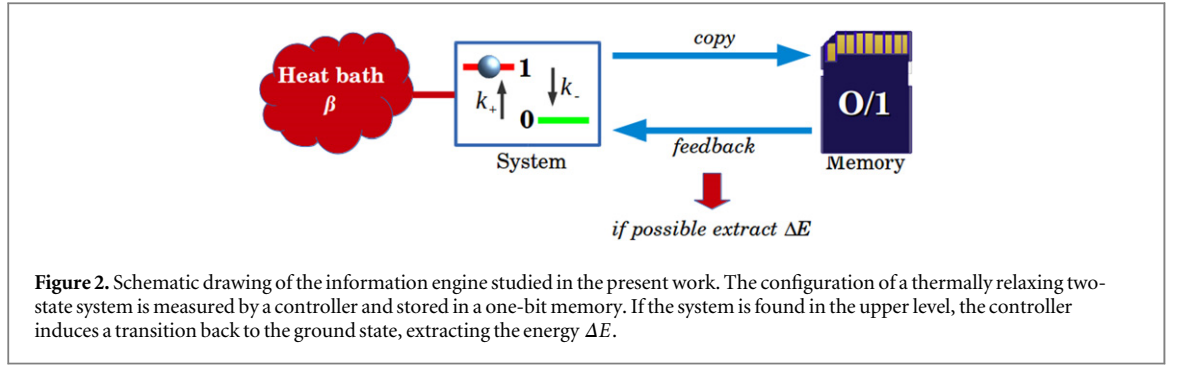
Subsequently this remarkable result was made more specific in various ways. For example, it was shown that one can construct feedback schemes that satisfy the inequality sharply [10, 11]. Moreover, the IFT was generalized to schemes with finite-time relaxation [12, 13] and continuous feedback schemes [14]. Very recently, the generalized IFT could also be confirmed experimentally [15].

The arising problem with these generalized Jarzynski equalities is that the tightness of the bound seems to depend on the specific feedback scheme, such as discrete and continuous feedback as well as memory tape models. This led Barato *et al* to look for a unifying master IFT [16–19]. To achieve this, they duplicated the configuration space, modeling bit flips of the memory by transitions between the two replicas. Doing so, they followed Smoluchowski's original idea of modeling the whole feedback loop as a physical device, defined as a stochastic Markov process.

In this paper, we follow these lines of thought, being interested in the *total* cost of operating an information engine with a finite memory as described above. Instead of duplicating the configuration space, we devise physically realizable stochastic processes of the joint system not only for relaxation, but also for the measurement process. We show that the generalized IFT for the relaxation of a system is always accompanied by an *opposite* IFT during the measurement carried out by the demon, restoring the second law for the joint system. This implies that a certain minimal amount of work has to be done on the memory, in accordance with Landauer's principle. The calculation of the total cost enables us to derive the efficiency of the information engine as a function of its cycle time. We also discuss the optimal cycle time and corresponding efficiency at which the extracted power is maximized.

## 2. Definition of a minimal model

In what follows we consider a system with two different energy levels separated by $\Delta E$. The system is coupled to a heat bath with inverse temperature $\beta$. The controller (demon) is implemented as a one bit memory. We devise a discrete feedback scheme evolving in three steps (see figure 2). First the system relaxes thermally without external influence. Then the actual energy level, denoted by 0 and 1, is copied to the memory of the controller. Finally, if the memory bit is 1, the controller induces a flip of the system $1 \rightarrow 0$, extracting the energy $\Delta E$; otherwise it does nothing.

**Figure 2.** Schematic drawing of the information engine studied in the present work. The configuration of a thermally relaxing two-state system is measured by a controller and stored in a one-bit memory. If the system is found in the upper level, the controller induces a transition back to the ground state, extracting the energy $\Delta E$.

Following Barato and Seifert [19], we are aiming to model these steps by stochastic Markov processes, including both the system and the memory. This defines a four-dimensional configuration space in which each step can be represented by a simple $4 \times 4$ matrix. In the following, we discuss each of the three steps in detail:

### 2.1. Relaxation
During relaxation, the system flips randomly according to the rules

$$0 \to 1 \quad \text{with rate } k_+$$
$$1 \to 0 \quad \text{with rate } k_-.$$

For convenience, we express these rates in terms of

$$k = k_+ + k_- \quad \text{and} \quad q = \frac{k_+}{k} \tag{2}$$

with $q < \frac{1}{2}$. After infinite time, the system eventually reaches an equilibrium state with the stationary probability distribution $P_{stat}^1 = 1 - P_{stat}^0 = q$. This means that the energy difference between the two levels is given by

$$\Delta E = \beta^{-1} \ln \frac{\bar{q}}{q} \qquad \left( \bar{q} = 1 - q \right). \tag{3}$$

Let us now describe the relaxation process in the composite configuration space of system and memory. Throughout this paper, we will use a canonical configuration basis ordered by

$$(\text{system state, memory bit}) = \{00, 01, 10, 11\}. \tag{4}$$

Since the memory is inactive during relaxation, the time evolution operator $\mathcal{L}_R$ for relaxation represented in this basis reads

$$\mathcal{L}_R = \underbrace{\begin{pmatrix} k_+ & -k_- \\ -k_+ & k_- \end{pmatrix}}_{\text{System}} \otimes \underbrace{\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}}_{\text{Memory}}$$

$$= \begin{pmatrix} k_+ & 0 & -k_- & 0 \\ 0 & k_+ & 0 & -k_- \\ -k_+ & 0 & k_- & 0 \\ 0 & -k_+ & 0 & k_- \end{pmatrix}. \tag{5}$$

After the relaxation time $t_R$, the corresponding transition matrix is given by

$$\mathcal{T}_R \left( t_R \right) = \exp \left( -\mathcal{L}_R t_R \right). \tag{6}$$

In the infinite-time relaxation limit $(t_R \to \infty)$, this transition matrix reduces to

$$\mathcal{T}_R = \lim_{t_R \to \infty} \mathcal{T}_R \left( t_R \right) = \begin{pmatrix} \bar{q} & 0 & \bar{q} & 0 \\ 0 & \bar{q} & 0 & \bar{q} \\ q & 0 & q & 0 \\ 0 & q & 0 & q \end{pmatrix}. \tag{7}$$

### 2.2. Measurement
A perfect measurement would faithfully copy the system state to the memory; i.e., if $m = 0, 1$ denotes the previous memory state and $s = 0, 1$ the actual system state, it would simply copy $(s, m) \mapsto (s, s)$. However, it is

well known that such a *perfect* measurement is irreversible, leading to a diverging entropy production [19, 20]. Therefore, one usually considers *imperfect* measurements

$$(s, m) \mapsto (s, s) \text{ with prob. } \bar{\epsilon} = 1 - \epsilon$$
$$(s, m) \mapsto (s, 1 - s) \text{ with prob. } \epsilon$$

with a small error probability $0 < \epsilon < 1/2$. Using the basis (4) this corresponds to the transition matrix

$$\mathcal{T}_M = \begin{pmatrix} \bar{\epsilon} & \bar{\epsilon} & 0 & 0 \\ \epsilon & \epsilon & 0 & 0 \\ 0 & 0 & \epsilon & \epsilon \\ 0 & 0 & \bar{\epsilon} & \bar{\epsilon} \end{pmatrix}. \tag{8}$$

Remarkably, this imperfect measurement process can be implemented by a stochastic Markov process as well, since $\mathcal{T}_M^2 = \mathcal{T}_M$. The corresponding time evolution operator reads

$$\mathcal{L}_M = k' \begin{pmatrix} \epsilon & -\bar{\epsilon} & 0 & 0 \\ -\epsilon & \bar{\epsilon} & 0 & 0 \\ 0 & 0 & \bar{\epsilon} & -\epsilon \\ 0 & 0 & -\bar{\epsilon} & \epsilon \end{pmatrix} \tag{9}$$

with a rate $k'$, and it is easy to show that the transition matrix $\mathcal{T}_M$ in equation (8) is retrieved in the limit of infinite measurement time:

$$\mathcal{T}_M = \lim_{t_M \to \infty} \mathcal{T}_M(t_M) = \lim_{t_M \to \infty} \exp\left(-\mathcal{L}_M t_M\right). \tag{10}$$

Thus, we succeeded in implementing the second step as a stochastic Markov process as well.

If the memory is considered as being in contact with some heat bath of inverse temperature $\beta$ during the stochastic measurement process, the time evolution defined above implies that the incorrectly measured state $(s, 1 - s)$ has a higher energy than the correctly measured state $(s, s)$ and that the corresponding energy difference between the two composite states is given by

$$\Delta E' = \beta^{-1} \ln \frac{\bar{\epsilon}}{\epsilon}. \tag{11}$$

### 2.3. Feedback

The purpose of the feedback is to use the information stored in the memory in order to extract energy from the system. If the preceding measurement was faithful, this would mean performing the transitions

$$00 \mapsto 00 \quad \text{without extraction of energy}$$
$$11 \mapsto 01 \quad \text{extracting the work } W_{\text{ex}} = \Delta E.$$

These transitions alone would be again irreversible, causing an infinite entropy production. However, if we add symmetric transitions in the (unlikely) case of erroneous measurements, namely,

$$10 \mapsto 10 \quad \text{without performing work}$$
$$01 \mapsto 11 \quad \text{performing work, i.e.,} W_{\text{ex}} = -\Delta E,$$
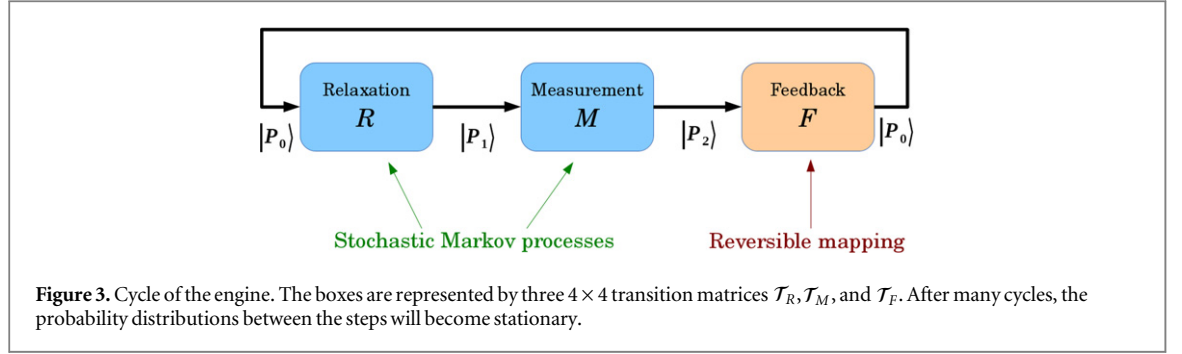
we obtain the feedback transition matrix

$$\mathcal{T}_F = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix}. \tag{12}$$

Thus the feedback process is carried out in such a way that the system state is flipped $(s \mapsto 1 - s)$ for $m = 1$, while it remains unchanged $(s \mapsto s)$ for $m = 0$. It is assumed that the feedback transition occurs instantaneously so that the total time $\tau$ of a complete cycle $\mathcal{T}_R \to \mathcal{T}_M \to \mathcal{T}_F$ is given by $\tau = t_R + t_M$.

Since $\mathcal{T}_F^2 = \mathbf{I}$, the feedback is fully reversible; hence, it does not produce entropy in the environment. Moreover, it is easy to see that it simply exchanges the second and the fourth component of a vector, and therefore it does not change the joint entropy of system and memory. However, as will be shown below, it generally changes the entropy of the subsystems.

Due to its reversible nature, the feedback as defined above cannot be implemented as a stochastic Markov process. However, we would like to point out that it is even possible to implement the feedback physically so that the entire chain of steps is represented cleanly as a sequence of stochastic processes. This can be done by replacing two subsequent cycles $\mathcal{T}_R \to \mathcal{T}_M \to \mathcal{T}_F \to \mathcal{T}_R \to \mathcal{T}_M \to \mathcal{T}_F \to \mathcal{T}_R$ equivalently by $\mathcal{T}_R \to \mathcal{T}_M \to \mathcal{T}_F \mathcal{T}_R \mathcal{T}_F \to \mathcal{T}_F \mathcal{T}_M \mathcal{T}_F \to \mathcal{T}_R$. Since $\mathcal{T}_{\tilde{R}} \equiv \mathcal{T}_F \mathcal{T}_R \mathcal{T}_F$ and $\mathcal{T}_{\tilde{M}} \equiv \mathcal{T}_F \mathcal{T}_M \mathcal{T}_F$ satisfy the stochasticity condition ($\mathcal{T}_{\tilde{R}}^2 = \mathcal{T}_{\tilde{R}}$, $\mathcal{T}_{\tilde{M}}^2 = \mathcal{T}_{\tilde{M}}$), the whole sequence of steps can be implemented by stochastic processes,

**Figure 3.** Cycle of the engine. The boxes are represented by three $4 \times 4$ transition matrices $\mathcal{T}_R, \mathcal{T}_M$, and $\mathcal{T}_F$. After many cycles, the probability distributions between the steps will become stationary.

providing a safe ground for the calculation of entropy production, mutual information, work, and heat. Having verified that this description is fully equivalent, we nevertheless keep the explicit feedback for simplicity in the original form.

Since the rates $k$ and $k'$ simply rescale $t_R$ and $t_M$, we will set

$$k = k' := 1 \tag{13}$$

throughout the paper. Thus, apart from $t_R$ and $t_M$, the model is controlled by only two parameters, namely, the relaxation parameter $q$ and the error probability $\epsilon$.

## 3. Stationary state

If the information engine runs repeatedly through many cycles, the probability distributions between the three steps will become stationary. The corresponding stationary probability distributions $|P_0\rangle, |P_1\rangle$, and $|P_2\rangle$ are represented as four-component vectors

$$|P_k\rangle = \left( P_k^{00}, P_k^{01}, P_k^{10}, P_k^{11} \right)^T, \quad (k = 0, 1, 2) \tag{14}$$

and are determined by the equations

$$\mathcal{T}_F \mathcal{T}_M \mathcal{T}_R |P_0\rangle = |P_0\rangle$$
$$\mathcal{T}_R \mathcal{T}_F \mathcal{T}_M |P_1\rangle = |P_1\rangle$$
$$\mathcal{T}_M \mathcal{T}_R \mathcal{T}_F |P_2\rangle = |P_2\rangle \tag{15}$$

with $|P_1\rangle = \mathcal{T}_R|P_0\rangle$ and $|P_2\rangle = \mathcal{T}_M|P_1\rangle$ (see figure 3).

The reduced stationary probability vectors of the system ($s$) and the memory ($m$) are given, respectively, as

$$|P_k^{(s)}\rangle = \left( P_k^{00} + P_k^{01}, P_k^{10} + P_k^{11} \right)^T$$
$$|P_k^{(m)}\rangle = \left( P_k^{00} + P_k^{10}, P_k^{01} + P_k^{11} \right)^T. \tag{16}$$

Clearly, the measurement does not modify the system state, meaning that $|P_1^{(s)}\rangle = |P_2^{(s)}\rangle$. Similarly, both the relaxation and feedback do not affect the memory; hence $|P_0^{(m)}\rangle = |P_1^{(m)}\rangle$ and $|P_2^{(m)}\rangle = |P_0^{(m)}\rangle$. As a result, in a stationary situation, the information acquired during the measurement is statistically the same in each cycle, implying that $|P_1^{(m)}\rangle = |P_2^{(m)}\rangle$.

As a simple example, first consider the case of infinite-time relaxation and measurement ($t_R, t_M \to \infty$), where the matrices $\mathcal{T}_R, \mathcal{T}_M$, and $\mathcal{T}_F$ are given by equations (7), (8), and (12). In this case the normalized stationary probability vectors turn out to be given by

$$|P_0\rangle = \left( \bar{\epsilon}\bar{q}, \quad \bar{\epsilon}q, \quad \epsilon q, \quad \epsilon\bar{q} \right)^T \tag{17}$$

$$|P_1\rangle = \left( \left( \bar{\epsilon}\bar{q} + \epsilon q \right)\bar{q}, \quad \left( \bar{\epsilon}q + \epsilon\bar{q} \right)\bar{q}, \quad \left( \bar{\epsilon}\bar{q} + \epsilon q \right)q, \quad \left( \bar{\epsilon}q + \epsilon\bar{q} \right)q \right)^T \tag{18}$$

$$|P_2\rangle = \left( \bar{\epsilon}\bar{q}, \quad \epsilon\bar{q}, \quad \epsilon q, \quad \bar{\epsilon}q \right)^T. \tag{19}$$

The corresponding reduced vectors for the system and the memory read

$$|P_0^{(s)}\rangle = \begin{pmatrix} \bar{\epsilon} \\ \epsilon \end{pmatrix}$$

$$|P_1^{(s)}\rangle = |P_2^{(s)}\rangle = \begin{pmatrix} \bar{q} \\ q \end{pmatrix}$$

$$|P_0^{(m)}\rangle = |P_1^{(m)}\rangle = |P_2^{(m)}\rangle = \begin{pmatrix} \bar{\epsilon}\bar{q} + \epsilon q \\ \bar{\epsilon}q + \epsilon\bar{q} \end{pmatrix}. \tag{20}$$

# 4. Entropy and entropy production

## 4.1. Shannon entropy

Given the stationary probability distributions $|P_k\rangle$, it is straightforward to compute the entropies of the system $(s)$, the memory $(m)$, and the joint system $(sm)$ between the steps in a stationary cycle, using the definition of the Shannon entropy

$$H = -\sum_c P^c \ln P^c, \tag{21}$$

where the sum runs over the vector components. Because of the aforementioned coincidence of various probability vectors we have $H_1^{(s)} = H_2^{(s)}$ and $H_1^{(m)} = H_2^{(m)} = H_3^{(m)}$. Furthermore, the reversibility of the feedback process guarantees that $H_2^{(sm)} = H_0^{(sm)}$.

For $t_R$, $t_M \to \infty$, the expressions for the Shannon entropies reduce to

$$H_0^{(s)} = h(\epsilon) + h(\bar{\epsilon})$$
$$H_{1,2}^{(s)} = h(q) + h(\bar{q}) \tag{22}$$

$$H_{0,1,2}^{(m)} = h(\epsilon\bar{q} + \bar{\epsilon}q) + h(\epsilon q + \bar{\epsilon}\bar{q}),$$
$$H_{0,2}^{(sm)} = h(q) + h(\bar{q}) + h(\epsilon) + h(\bar{\epsilon})$$
$$H_1^{(sm)} = h(q) + h(\bar{q}) + h(\epsilon\bar{q} + \bar{\epsilon}q) + h(\epsilon q + \bar{\epsilon}\bar{q}), \tag{23}$$

where we used the notation $h(p) := -p \ln p$.

During the relaxation process, where the system tries to restore the equilibrium distribution from the overpopulated ground state after energy extraction, the system entropy $H^{(s)}$ is expected to increase, provided that the error probability $\epsilon$ is sufficiently small ($\epsilon < q$). The same applies to the composite entropy $H^{(sm)}$.

To summarize, the entropy changes during relaxation (R), measurement (M), and feedback (F) are given by

$$\Delta H_R^{(s)} > 0, \ \Delta H_M^{(s)} = 0, \ \Delta H_F^{(s)} = -\Delta H_R^{(s)} < 0,$$
$$\Delta H_R^{(m)} = \Delta H_M^{(m)} = \Delta H_F^{(m)} = 0,$$
$$\Delta H_R^{(sm)} > 0, \ \Delta H_M^{(sm)} = -\Delta H_R^{(sm)}, \ \Delta H_F^{(sm)} = 0. \tag{24}$$

## 4.2. Mutual information

With these expressions, it is straightforward to compute the mutual information

$$I_k = H_k^{(s)} + H_k^{(m)} - H_k^{(sm)} \geqslant 0, \quad (k = 0, 1, 2) \tag{25}$$

which is a measure for the correlation between system and memory. This correlation is expected to build up during the measurement and then to decrease during feedback and relaxation, implying the inequalities

$$\Delta I_R < 0, \quad \Delta I_M > 0, \quad \Delta I_F < 0. \tag{26}$$

It is interesting to note that the change of the composite entropy is purely given by the amount of mutual information acquired during the measurement, i.e.,

$$\Delta H_M^{(sm)} = -\Delta I_M \quad \text{and} \quad \Delta H_R^{(sm)} = \Delta I_M. \tag{27}$$

For $t_R$, $t_M \to \infty$, we have

$$
\begin{aligned}
I_0 &= h\left(\epsilon\bar{q} + \bar{\epsilon}q\right) + h\left(\epsilon q + \bar{\epsilon}\bar{q}\right) - h(q) - h\left(\bar{q}\right) \\
I_1 &= 0 \\
I_2 &= h\left(\epsilon\bar{q} + \bar{\epsilon}q\right) + h\left(\epsilon q + \bar{\epsilon}\bar{q}\right) - h(\epsilon) - h\left(\bar{\epsilon}\right).
\end{aligned}
\tag{28}
$$

Note that the result $I_1 = 0$ is true only if the relaxation time is infinite, which obviously destroys all correlations between system and memory.

### 4.3. Entropy production

Let us now turn to entropy production. According to Schnakenberg [21–23], whenever the system or the memory jumps spontaneously from the configuration $c$ to another configuration $c'$, the amount of entropy

$$
\Delta H^{\mathrm{env}}_{c \to c'}(t) = \ln \frac{w_{c \to c'}(t)}{w_{c' \to c}(t)}
\tag{29}
$$

is generated in the environment. Here, $w_{c \to c'}(t)$ denotes the transition rate at time $t$. Therefore, the mean entropy production rate is given by

$$
\frac{\mathrm{d}}{\mathrm{d}t} H^{\mathrm{env}} = \sum_{c \neq c'} P^c(t) w_{c \to c'}(t) \ln \frac{w_{c \to c'}(t)}{w_{c' \to c}(t)}.
\tag{30}
$$

In an arbitrary nonequilibrium system, one would have to solve the master equation, plug the solution into the equation above, and integrate the resulting expression over a certain window of time. However, in the present case this is not necessary, since the model is so simple that the rates happen to obey *detailed balance*, defined as $P^c_{stat} w_{c \to c'} = P^{c'}_{stat} w_{c' \to c}$ in the stationary equilibrium state. In this case, it is therefore straightforward to rewrite the equation given above as

$$
\begin{aligned}
\frac{\mathrm{d}}{\mathrm{d}t} H^{\mathrm{env}} &= \sum_{c \neq c'} \left( P^{c'} w_{c' \to c} - P^c w_{c \to c'} \right) \ln P^c_{stat} \\
&= \sum_c \frac{\mathrm{d}P^c}{\mathrm{d}t} \ln P^c_{stat},
\end{aligned}
\tag{31}
$$

allowing us to compute the average entropy production in each step directly without integration by means of

$$
\Delta H^{\mathrm{env}} = \sum_c \left( P^c_{final} - P^c_{init} \right) \ln P^c_{stat}.
\tag{32}
$$

Here $P^c_{init}$ and $P^c_{final}$ denote the initial and the final probabilities for a finite time span, while $P^c_{stat}$ is the stationary probability distribution that would emerge after an infinite long time. For example, during relaxation, $P^c_{stat}$ can be obtained by taking $t_R \to \infty$ in the expression for $P^c_1$ in equation (18). Similarly, during measurement, $P^c_{stat} = \lim_{t_M \to \infty} P^c_2$ with $P^c_2$ given in equation (19).

With the above formula, we obtain the following expressions for the entropy production in each process:
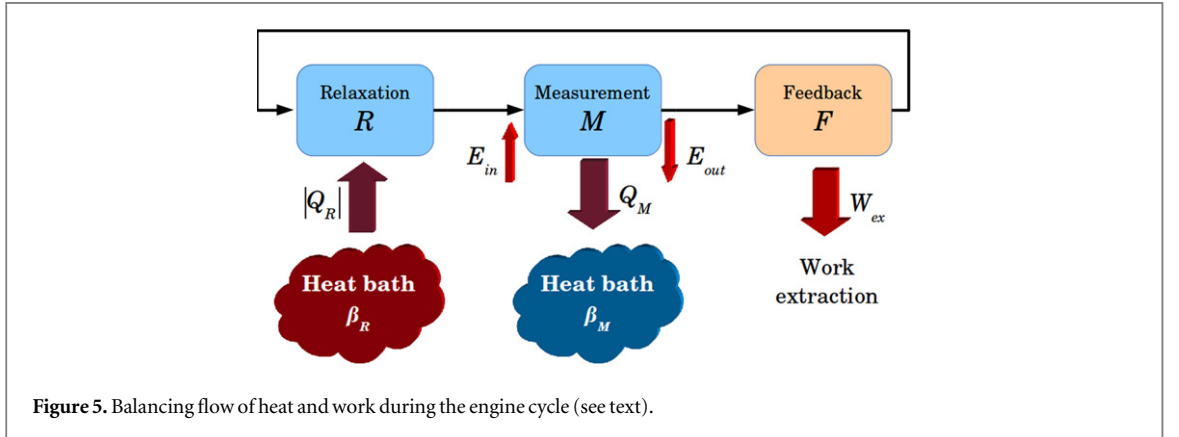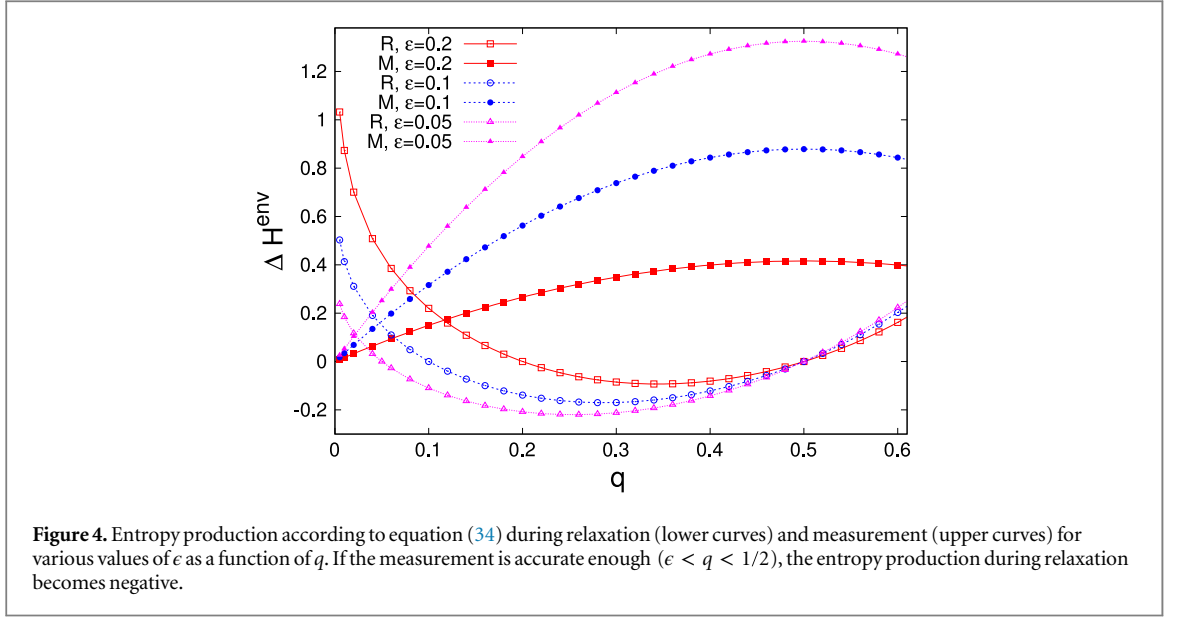
$$
\begin{aligned}
\Delta H^{\mathrm{env}}_R &= \left( P^{10}_0 + P^{11}_0 - P^{10}_1 - P^{11}_1 \right) \ln \frac{\bar{q}}{q} \\
\Delta H^{\mathrm{env}}_M &= \left( P^{01}_1 + P^{10}_1 - P^{01}_2 - P^{10}_2 \right) \ln \frac{\bar{\epsilon}}{\epsilon} \\
\Delta H^{\mathrm{env}}_F &= 0.
\end{aligned}
\tag{33}
$$

Note that these expressions hold for any finite $t_R$ and $t_M$. The last result is obvious, since the feedback with $\mathcal{T}^2_F = \mathbf{I}$ is a reversible operation.

For $t_R$, $t_M \to \infty$, by inserting equations (17)–(19) we get explicit expressions

$$
\begin{aligned}
\Delta H^{\mathrm{env}}_R &= -(q - \epsilon) \ln \frac{\bar{q}}{q} \\
\Delta H^{\mathrm{env}}_M &= 2q\bar{q}\left( \bar{\epsilon} - \epsilon \right) \ln \frac{\bar{\epsilon}}{\epsilon},
\end{aligned}
\tag{34}
$$

which are plotted for various error probabilities in figure 4. As one can see, for $\epsilon < q < 1/2$ the entropy production during relaxation $\Delta H^{\mathrm{env}}_R$ is negative, meaning that the engine imports entropy (heat) from the environment rather than producing it. Obviously, this is the regime of interest where we would like to operate our information engine. However, as can be seen, the negative entropy production during relaxation is always overcompensated by a positive one during the measurement, which is consistent with the second law of thermodynamics. Notice that the entropy production during measurement is always positive, since $\epsilon < 1/2$.

**Figure 4.** Entropy production according to equation (34) during relaxation (lower curves) and measurement (upper curves) for various values of $\epsilon$ as a function of $q$. If the measurement is accurate enough ($\epsilon < q < 1/2$), the entropy production during relaxation becomes negative.



**Figure 5.** Balancing flow of heat and work during the engine cycle (see text).

## 5. Work extraction and supply

By virtue of Clausius' law $dQ = T\,dH$, the produced entropy can be translated directly into an amount of heat. In most studies, it is usually assumed that the temperatures of the system and the memory are identical. However, for the sake of generality, let us allow the temperatures to be different, assigning $\beta_R = 1/T_R$ during relaxation and $\beta_M = 1/T_M$ during measurement, as sketched in figure 5, including the conventional setup of information engines with a single reservoir as a special case. Thus the respective heat contributions averaged over all possible stochastic trajectories are given by

$$\langle Q_R \rangle = \beta_R^{-1} \Delta H_R^{\text{env}}, \qquad \langle Q_M \rangle = \beta_M^{-1} \Delta H_M^{\text{env}}. \tag{35}$$

Here we use the sign convention that heat flowing away from the engine into the environment has a positive sign; i.e., we expect $\langle Q_R \rangle$ to be negative and $\langle Q_M \rangle$ to be positive.

In order to maintain stationarity of the system after one engine cycle, the average work $\langle W_{\text{ex}} \rangle$ extracted during feedback should exactly balance the average heat $\langle Q_R \rangle$ during relaxation; i.e.,

$$\langle W_{\text{ex}} \rangle = -\langle Q_R \rangle > 0. \tag{36}$$

Consequently, the measurement process does not change the energy of the system. This requires an additional influx of energy in the form of extra work $\langle W_{sup} \rangle$ into the memory, which is necessary to compensate the loss of heat $\langle Q_M \rangle$ flowing away to the environment during the measurement process:

$$\langle W_{sup} \rangle = \langle Q_M \rangle > 0. \tag{37}$$

The average net work performed by the machine, defined as the difference of extracted and supplied work, is therefore given by

$$\langle W_{\text{net}} \rangle = \langle W_{\text{ex}} \rangle - \langle W_{sup} \rangle = -\langle Q_R \rangle - \langle Q_M \rangle. \tag{38}$$

Note that the net work can change its sign, depending on the choice of the parameters $q$, $\epsilon$, $t_R$, and $t_M$.

Let us first compute the extracted work. According to section 2.3, $\langle W_{\text{ex}} \rangle$ is given by $(P_2^{11} - P_2^{01})\Delta E$, where $\Delta E = \beta_R^{-1} \ln \frac{\bar{q}}{q}$ (see (3)). Using $|P_1^{(s)}\rangle = |P_2^{(s)}\rangle$ (no system change during measurement), together with the feedback identities $P_2^{01} = P_0^{11}$ and $P_2^{10} = P_0^{10}$ (flip $s$ only when $m = 1$), it is easy to show explicitly that

$$\langle W_{\text{ex}} \rangle = -\beta_R^{-1} \Delta H_R^{\text{env}} = -\langle Q_R \rangle. \tag{39}$$

The additional work $\langle W_{sup} \rangle$ supplied during the measurement process can be interpreted as the energy needed to operate the measurement device. Technically this contribution comes from the fact that the energy levels of the joint system are different during relaxation and measurement so that extra energy is needed to move them around. For example, this could be done by applying an external potential in order to make the energy level of the erroneous composite state $(s, 1 - s)$ higher than that of the correctly measured state $(s,s)$ by the amount of $\Delta E' = \beta_M^{-1} \ln \frac{\bar{\epsilon}}{\epsilon}$. When the external potential is turned on just before the measurement, the average energy of the composite of system and memory increases by $\langle E_{in} \rangle$, and, similarly, it loses the energy $\langle E_{out} \rangle$ when the potential is turned off at the end of the measurement:

$$\langle E_{in} \rangle = \frac{P_1^{01} + P_1^{10}}{\beta_M} \ln \frac{\bar{\epsilon}}{\epsilon}, \quad \langle E_{out} \rangle = \frac{P_2^{01} + P_2^{10}}{\beta_M} \ln \frac{\bar{\epsilon}}{\epsilon}. \tag{40}$$

Comparing the difference $\langle W_{sup} \rangle = \langle E_{in} \rangle - \langle E_{out} \rangle$ with equation (33), one can see immediately that

$$\langle W_{sup} \rangle = \beta_M^{-1} \Delta H_M^{\text{env}} = \langle Q_M \rangle. \tag{41}$$

In the limit $t_R$, $t_M \to \infty$, the average work contributions read

$$\langle W_{\text{ex}} \rangle = \frac{q - \epsilon}{\beta_R} \ln \frac{\bar{q}}{q}$$

$$\langle W_{sup} \rangle = \frac{2q\bar{q}\left(\bar{\epsilon} - \epsilon\right)}{\beta_M} \ln \frac{\bar{\epsilon}}{\epsilon}. \tag{42}$$

## 6. Thermodynamic second laws and fluctuation theorems

The total entropy production of the whole setup during relaxation (R) and measurement (M) is given by

$$\Delta H_R^{\text{tot}} = \Delta H_R^{(sm)} + \Delta H_R^{\text{env}} \tag{43}$$

$$\Delta H_M^{\text{tot}} = \Delta H_M^{(sm)} + \Delta H_M^{\text{env}}. \tag{44}$$

According to equation (27), the entropy differences in the joint system are given solely in terms of the mutual information difference

$$\Delta H_R^{sm} = -\Delta H_M^{sm} = \Delta I_M \tag{45}$$

while the entropy differences in the environment are given in terms of the transferred heat by $\Delta H_{R,M}^{\text{env}} = \beta_{R,M} \langle Q_{R,M} \rangle$. Using equations (39) and (41), we arrive at

$$\Delta H_R^{\text{tot}} = +\Delta I_M - \beta_R \langle W_{\text{ex}} \rangle$$

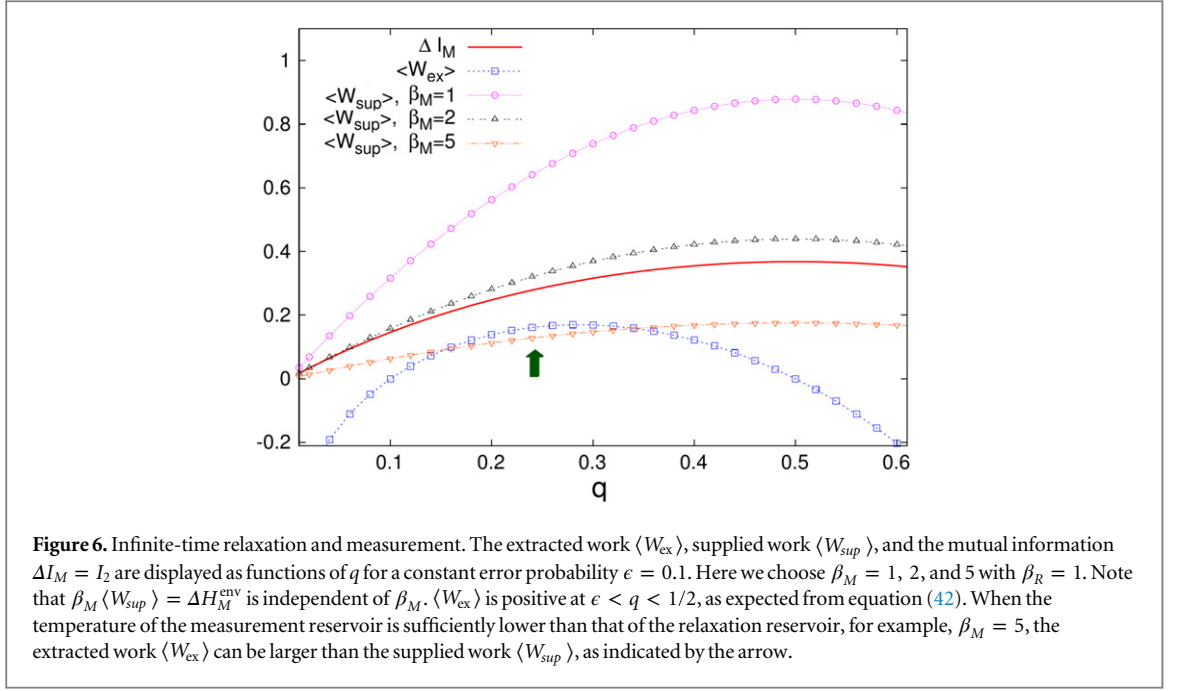$$\Delta H_M^{\text{tot}} = -\Delta I_M + \beta_M \langle W_{sup} \rangle. \tag{46}$$

For the total system, including the environment, the second law of thermodynamics should be satisfied for each process; i.e.,

$$\Delta H_R^{\text{tot}} \geqslant 0 \quad \text{and} \quad \Delta H_M^{\text{tot}} \geqslant 0, \tag{47}$$

or, equivalently,

$$\beta_R \langle W_{\text{ex}} \rangle \leqslant \Delta I_M \leqslant \beta_M \langle W_{sup} \rangle. \tag{48}$$

Most existing studies are only interested in the first inequality during relaxation, while the other one during measurement is ignored. The purpose of this work is to point out that there is a *second* inequality for the measurement process as well, and that the two inequalities are complementary with respect to each other. More specifically, if $\beta_R = \beta_M$, the extracted work is bounded from above by the mutual information, while the work required to operate the memory is bounded from below by the same threshold. This means that the setup cannot be used to gain work out of nothing, $\langle W_{\text{net}} \rangle \leqslant 0$, as expected by the second law. However, if the two reservoir temperatures were different $(\beta_R < \beta_M)$, the extracted work could be larger than the supplied one, in which case the system operates like a conventional heat engine (see figure 6).

**Figure 6.** Infinite-time relaxation and measurement. The extracted work $\langle W_{ex} \rangle$, supplied work $\langle W_{sup} \rangle$, and the mutual information $\Delta I_M = I_2$ are displayed as functions of $q$ for a constant error probability $\epsilon = 0.1$. Here we choose $\beta_M = 1$, 2, and 5 with $\beta_R = 1$. Note that $\beta_M \langle W_{sup} \rangle = \Delta H_M^{env}$ is independent of $\beta_M$. $\langle W_{ex} \rangle$ is positive at $\epsilon < q < 1/2$, as expected from equation (42). When the temperature of the measurement reservoir is sufficiently lower than that of the relaxation reservoir, for example, $\beta_M = 5$, the extracted work $\langle W_{ex} \rangle$ can be larger than the supplied work $\langle W_{sup} \rangle$, as indicated by the arrow.

It is almost trivial to construct the IFTs through the standard approach of stochastic thermodynamics [23] by considering the heat along all possible trajectories in the composite configurational state space. With an appropriate definition of the Shannon entropy for a given trajectory [23], one can easily get the fluctuation theorems for the total entropy production for each process as

$$\langle e^{-\Delta H_R^{tot}(\text{traj.})} \rangle = 1, \qquad \langle e^{-\Delta H_M^{tot}(\text{traj.})} \rangle = 1. \tag{49}$$

The second laws in equation (47) are simple consequences of the IFTs with $\Delta H_{R,M}^{tot} = \langle \Delta H_{R,M}^{tot}(\text{traj.}) \rangle$. It is rather tricky to find the IFTs in terms of work and mutual information, because it requires an equilibrium state as an initial condition. This is the case only when $t_R$ becomes infinite so that the system is in equilibrium at the start of the measurement as well as at the beginning of the feedback. However, note that the bounds for works in equation (48) are valid, even if $t_R$ is finite.

## 7. Finite-time relaxation and measurement

In practice, an engine is only useful if the cycle time $\tau$ is finite. Thus, it is obviously of interest to derive all physical quantities as a function of the cycle time. This allows one to find the optimum for maximal power generation, as will be discussed in the next section.

It is straightforward to obtain the transition matrices for relaxation and measurement for finite time spans $t_R$ and $t_M$:

$$\mathcal{T}_R(t_R) = \begin{pmatrix} \mathcal{R} + \bar{\mathcal{R}}\bar{q} & 0 & \bar{\mathcal{R}}\bar{q} & 0 \\ 0 & \mathcal{R} + \bar{\mathcal{R}}\bar{q} & 0 & \bar{\mathcal{R}}\bar{q} \\ \bar{\mathcal{R}}q & 0 & \mathcal{R} + \bar{\mathcal{R}}q & 0 \\ 0 & \bar{\mathcal{R}}q & 0 & \mathcal{R} + \bar{\mathcal{R}}q \end{pmatrix}$$
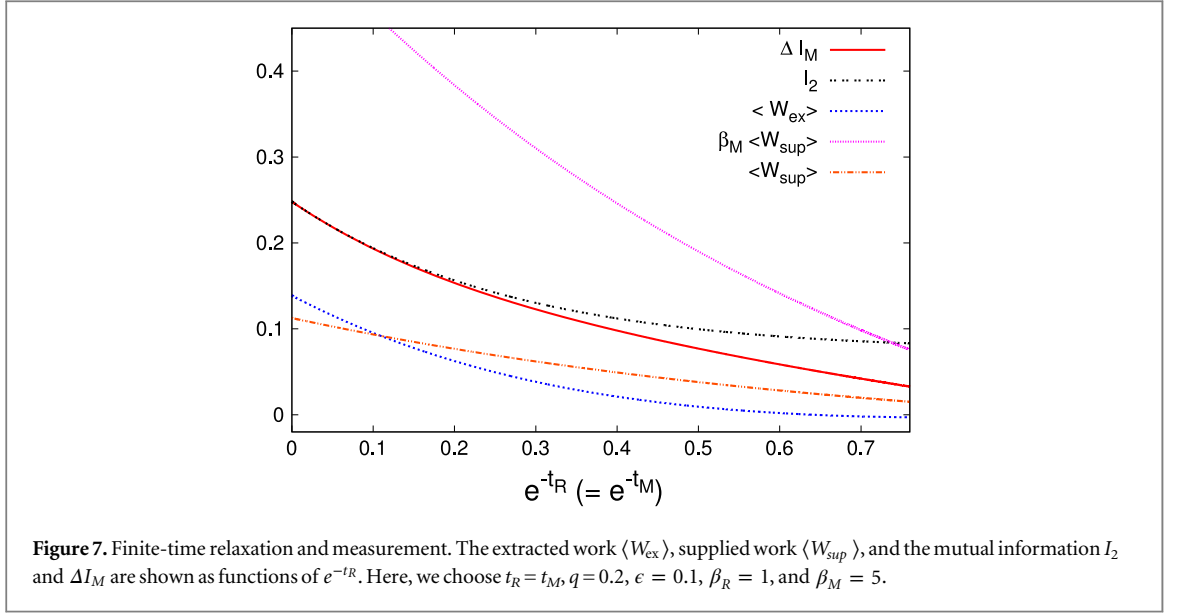
$$\mathcal{T}_M(t_M) = \begin{pmatrix} \mathcal{M} + \bar{\mathcal{M}}\bar{\epsilon} & \bar{\mathcal{M}}\bar{\epsilon} & 0 & 0 \\ \bar{\mathcal{M}}\epsilon & \mathcal{M} + \bar{\mathcal{M}}\epsilon & 0 & 0 \\ 0 & 0 & \mathcal{M} + \bar{\mathcal{M}}\epsilon & \bar{\mathcal{M}}\epsilon \\ 0 & 0 & \bar{\mathcal{M}}\bar{\epsilon} & \mathcal{M} + \bar{\mathcal{M}}\bar{\epsilon} \end{pmatrix}, \tag{50}$$

where

$$\mathcal{R} := e^{-t_R} \quad \text{and} \quad \bar{\mathcal{R}} := 1 - \mathcal{R} \tag{51}$$

$$\mathcal{M} := e^{-t_M} \quad \text{and} \quad \bar{\mathcal{M}} := 1 - \mathcal{M}. \tag{52}$$

Note that the transition probability from $(s, 1 - s)$ to $(s, s)$ during measurement, $\bar{\mathcal{M}}\bar{\epsilon}$, is smaller than $\bar{\epsilon}$, which means that the measurement for finite $t_M$ is less accurate than in the limit of infinite time.

**Figure 7.** Finite-time relaxation and measurement. The extracted work $\langle W_{ex} \rangle$, supplied work $\langle W_{sup} \rangle$, and the mutual information $I_2$ and $\Delta I_M$ are shown as functions of $e^{-t_R}$. Here, we choose $t_R = t_M$, $q = 0.2$, $\epsilon = 0.1$, $\beta_R = 1$, and $\beta_M = 5$.

Solving equation (15), one can find explicit but complicated expressions for all stationary distributions such as $|P_0\rangle$, $|P_1\rangle$, and $|P_2\rangle$ for finite $t_R$ and $t_M$, which are not shown here explicitly. The heat dissipation during the finite-time relaxation and measurement can be obtained from equation (33), while the extracted work and the supplied work are given by equations (36) and (37).

Using the relation $|P_2\rangle = \mathcal{T}_M |P_1\rangle$ we find that

$$\langle W_{sup} \rangle = \frac{\mathcal{E} - \epsilon}{\beta_M} \bar{\mathcal{M}} \ln \frac{\bar{\epsilon}}{\epsilon}. \tag{53}$$

Here $\mathcal{E} \equiv P_1^{01} + P_1^{10}$ is explicitly given by

$$\mathcal{E}\left(t_R, t_M\right) = \frac{\bar{\mathcal{R}}\left(\bar{\epsilon}q + \epsilon\bar{q}\right) + \alpha\left(\bar{\mathcal{R}}q + \mathcal{R}\bar{\mathcal{M}}\epsilon\right)}{1 - \alpha\mathcal{R}\mathcal{M}}, \tag{54}$$

with $\alpha = \mathcal{R} + \bar{\mathcal{R}}\left(\bar{q} - q\right)\left(\bar{\epsilon} - \epsilon\right)$. In a similar manner, we obtain $\langle W_{ex} \rangle$ as the function of $\mathcal{E}$:

$$\langle W_{ex} \rangle = \frac{q - \epsilon - \mathcal{M}(\mathcal{E} - \epsilon)}{\beta_R} \bar{\mathcal{R}} \ln \frac{\bar{q}}{q}. \tag{55}$$

In the limit of $t_R, t_M \to \infty$ ($\mathcal{R}, \mathcal{M} \to 0$), we consistently recover equation (42).
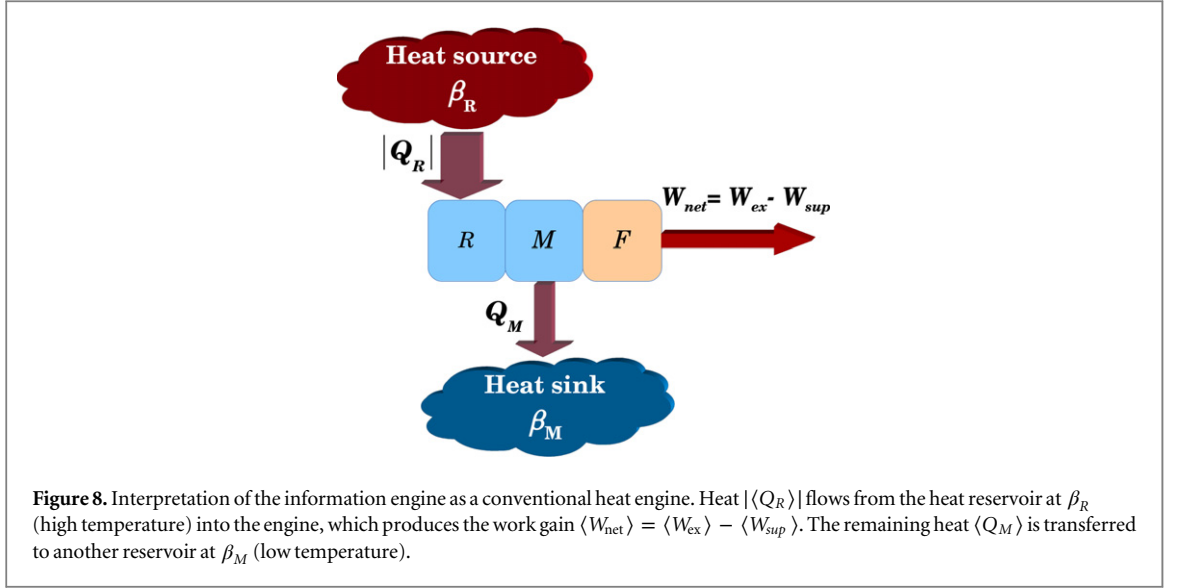
Note that the finite-time works in (53) and (55) decrease monotonously with $\mathcal{R}$ and $\mathcal{M}$ (see figure 7). Moreover, since the correlation between system and memory builds up continuously during the measurement process, it is obvious that $\Delta I_M$ decreases with $\mathcal{M}$, remaining positive by definition. The positivity of $\Delta I_M$ guarantees that $\langle W_{sup} \rangle$ is also positive. On the other hand, $\langle W_{ex} \rangle$ can be negative for short-time measurement and relaxation, as its upper bound $\Delta I_M$ approaches zero for $t_M \to 0$.

The monotonous dependence shown in figure 7 suggests that $\langle W_{ex} \rangle$ becomes maximal for maximal measurement accuracy ($t_M \to \infty$) and full relaxation ($t_R \to \infty$) in order to redistribute and pump the overpopulated ground state $s = 0$ back to the energetically excited state $s = 1$. Therefore, both limits $t_M, t_R \to \infty$ have to be carried out simultaneously. To establish this combined limit conveniently, let us from now on set

$$t_R = t_M = \tau/2, \tag{56}$$

meaning that $\mathcal{R} = \mathcal{M} = e^{-\tau/2}$. With this convention we expect $\langle W_{ex} \rangle$ to be maximal in the limit of infinite cycle time ($\tau \to \infty$). Moreover, as $\tau$ decreases, we expect $\langle W_{ex} \rangle$ to decrease and eventually to become negative.

If $\langle W_{ex} \rangle > \langle W_{sup} \rangle$, the system operates like a conventional heat engine. For infinite $\tau$ the net work $\langle W_{net} \rangle$ is maximal, but the power (net work per unit time) vanishes. For finite but sufficiently large $\tau$ and properly chosen parameters, the system still produces a positive net work; hence, the power is positive. However, as can be seen in figure 7, the curves for $\langle W_{ex} \rangle$ and $\langle W_{sup} \rangle$ cross each other at some finite cycle time $\tau = \tau_s$. At this point we no longer obtain any net work from the engine; hence, the power vanishes again. Consequently, there will be a particular cycle time in between, at which the power of the engine is maximal. In the next section, we will discuss this aspect in more detail.

**Figure 8.** Interpretation of the information engine as a conventional heat engine. Heat $|\langle Q_R \rangle|$ flows from the heat reservoir at $\beta_R$ (high temperature) into the engine, which produces the work gain $\langle W_{\text{net}} \rangle = \langle W_{\text{ex}} \rangle - \langle W_{sup} \rangle$. The remaining heat $\langle Q_M \rangle$ is transferred to another reservoir at $\beta_M$ (low temperature).

## 8. Efficiency

Let us now assume that the information engine operates in a regime where the net work is positive. In this case the whole setup can be interpreted as a conventional heat engine, as sketched in figure 8. As $\beta_R < \beta_M$, the upper reservoir for the relaxation process plays the role of a high-temperature heat source, while the lower reservoir in contact with the memory device acts as a heat sink. The *efficiency* of this heat engine in a single cycle is defined in the usual way as

$$\eta(\tau) = \frac{\langle W_{\text{net}} \rangle}{|\langle Q_R \rangle|} = 1 - \frac{\langle W_{sup} \rangle}{\langle W_{\text{ex}} \rangle}. \tag{57}$$

Using equations (53) and (55) the efficiency can be rewritten as

$$\eta(\tau) = 1 - \frac{\beta_R}{\beta_M} \lambda(\tau), \tag{58}$$

where

$$\lambda(\tau) = \frac{(\mathcal{E} - \epsilon)\bar{M} \ln\left(\bar{\epsilon}/\epsilon\right)}{[q - \epsilon - \mathcal{M}(\mathcal{E} - \epsilon)]\bar{R} \ln\left(\bar{q}/q\right)}. \tag{59}$$
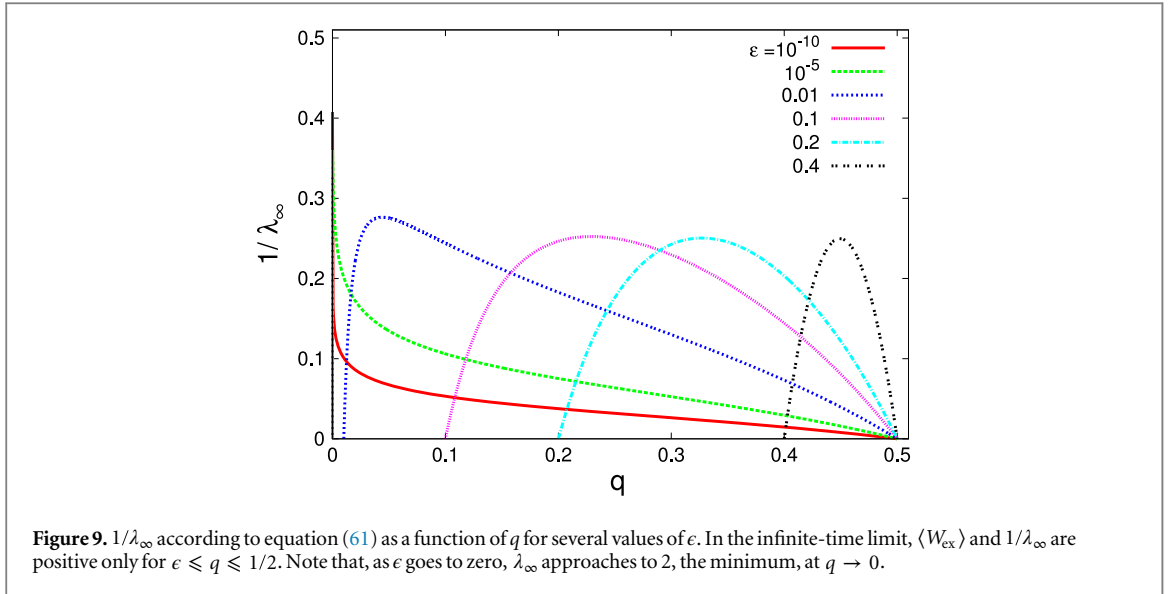
Note that the relaxation and measurement processes are not quasi-static, so that even in the limit $\tau \to \infty$ the engine never reaches the Carnot efficiency $\eta_c = 1 - \frac{\beta_R}{\beta_M}$. Instead, we find that the efficiency is limited by a different upper bound $\eta_{\text{max}}$, which can be computed as follows. According to the monotonicity arguments discussed in the preceding section, $\eta(\tau)$ is expected to become maximal in the limit $\tau \to \infty$. This suggests that

$$\eta(\tau) \leqslant \eta_{\text{max}} \equiv \lim_{\tau \to \infty} \eta(\tau) = 1 - \frac{\beta_R}{\beta_M} \lambda_\infty, \tag{60}$$

where

$$\lambda_\infty = \lim_{\tau \to \infty} \lambda(\tau) = \frac{2q\bar{q}\left(\bar{\epsilon} - \epsilon\right)\ln\left(\bar{\epsilon}/\epsilon\right)}{(q - \epsilon)\ln\left(\bar{q}/q\right)}. \tag{61}$$

Figure 9 shows $1/\lambda_\infty$ as a function of $q$ for several values of $\epsilon$. It is obvious that $1/\lambda_\infty$ is positive for $\epsilon \leqslant q \leqslant 1/2$, with a unimodal shape due to the similar behavior of $\langle W_{\text{ex}} \rangle$ demonstrated in figure 6. As $\epsilon \to 0$, $1/\lambda_\infty$ approaches zero, except for $q \approx 0$. In the limit of both $q \to 0$ and $\epsilon \to 0$, $\lambda_\infty$ approaches a constant bounded from below by $\lambda_\infty \geqslant 2$. Consequently, in order to obtain a positive work gain, the difference of temperatures should be at least $\beta_M/\beta_R \geqslant 2$ for the infinite-time process. In short, we find that the efficiency of the information engine is bounded by

**Figure 9.** $1/\lambda_\infty$ according to equation (61) as a function of $q$ for several values of $\epsilon$. In the infinite-time limit, $\langle W_{\mathrm{ex}} \rangle$ and $1/\lambda_\infty$ are positive only for $\epsilon \leqslant q \leqslant 1/2$. Note that, as $\epsilon$ goes to zero, $\lambda_\infty$ approaches to 2, the minimum, at $q \to 0$.

$$\eta_{\max} \leqslant 1 - \frac{2\beta_R}{\beta_M} < \eta_c. \tag{62}$$

In conventional heat engines operating with a finite cycle time, thermodynamic processes are no longer quasi-static, leading to an efficiency below the Carnot bound $\eta_c$. In our case, we also find that $\eta(\tau)$ becomes maximal in the limit $\tau \to \infty$. However, in contrast to the Carnot limit, where the entropy production vanishes, the entropy production per cycle in our model $\Delta H_R^{\mathrm{tot}} + \Delta H_M^{\mathrm{tot}}$ becomes *also maximal* at $\tau = \infty$. This is one of the crucial features of our information engine, which is totally distinct from conventional heat engines. Nevertheless, the entropy production *rate* (per unit time) decreases with increasing $\tau$ and finally vanishes at $\tau = \infty$. Hence, one may also say that the maximum efficiency is found at the minimum entropy production rate.

Similarly, the average power gain, $\langle P \rangle \equiv \langle W_{\mathrm{net}} \rangle / \tau$, in fact vanishes for $\tau \to \infty$, because $\langle W_{\mathrm{net}} \rangle$ remains finite in this limit. For a realistic engine, we usually want to optimize the power gain, trading off the efficiency against the cycle time. As expected, $\langle P \rangle$ is maximized at some *finite* time, $\tau_{op}$ between $\tau_s$ and $\infty$, as shown in figure 10(a).

The efficiency of heat engines at maximal power has been studied previously in [24–26]. Especially, for the Curzon–Ahlborn (CA) endoreversible model [24], it is well–known that the efficiency $\eta_{\mathrm{CA}}$ at the optimal power is given by $\eta_{\mathrm{CA}} = 1 - \sqrt{1 - \eta_c}$. In figure 11, we plot the efficiency at optimal power $\eta_{op} = \eta(\tau = \tau_{op})$, according to equation (58), as a function of $\eta_{\max}$ instead of $\eta_c$. It turns out that the functional behavior of $\eta_{op}$ is completely different from $\eta_{\mathrm{CA}}$.

In more general situations, it has been found that the efficiency at the maximum power obeys a universal form, $\eta_{op} = \eta_c/2 + O(\eta_c^2)$ for small $\eta_c$, when the engine and heat baths are strongly coupled [25, 26]. Our information engine exhibits a completely distinct behavior, even for small $\eta_{\max}$. As seen in figure 11, $\eta_{op}$ is more or less the same as $\eta_{\max}$ in this regime.

In order to investigate this unusual behavior in more detail, we now examine $\eta_{op}$ analytically for small $\eta_{\max}$. As the stall time $\tau_s$, and thus $\tau_{op}$ $(> \tau_s)$ becomes large, a small $\mathcal{R}$ expansion may be valid for small $\eta_{\max}$. Expanding $\eta(\tau)$ in equation (58) up to the linear order of $\mathcal{R} = e^{-\tau/2}$, we get

$$\eta(\tau) \approx \eta_{\max} - \left( \frac{\mathcal{E}_1}{\mathcal{E}_0 - \epsilon} + \frac{\mathcal{E}_0 - \epsilon}{q - \epsilon} \right) e^{-\tau/2}, \tag{63}$$
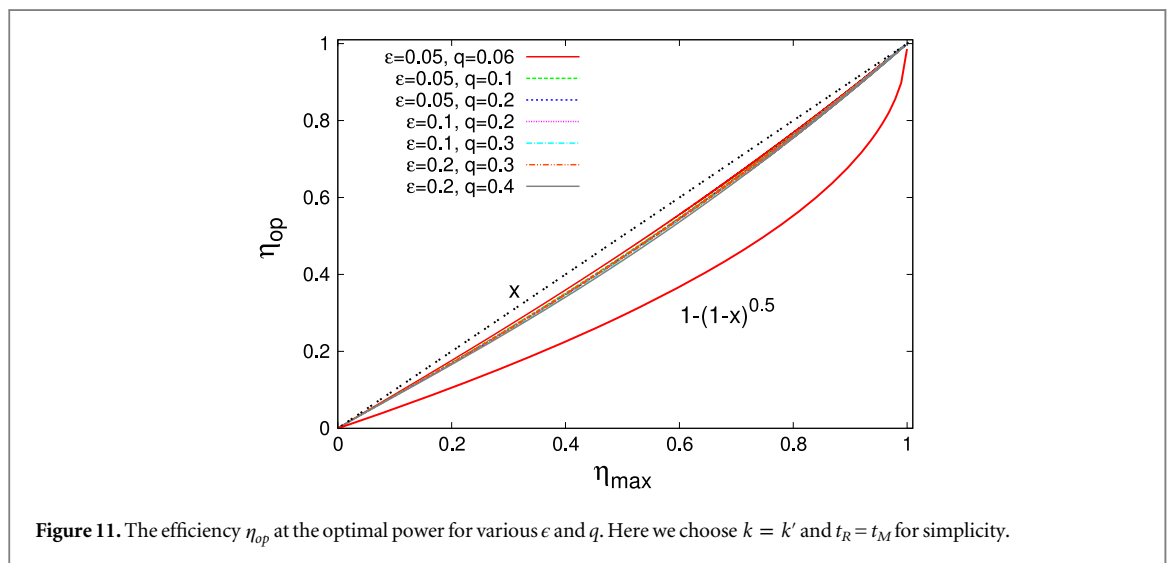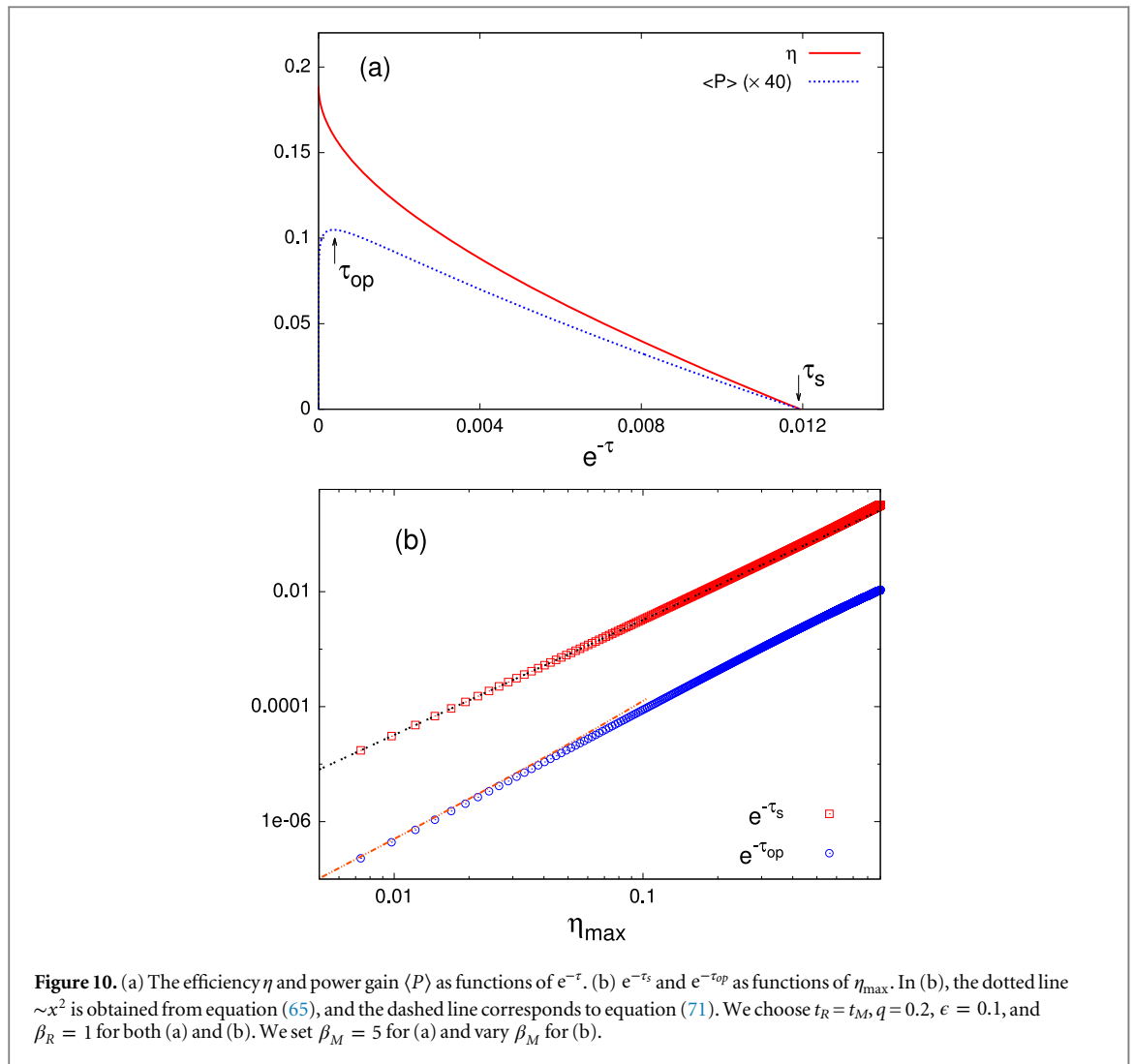
where we have used $\mathcal{E} \approx \mathcal{E}_0 + \mathcal{E}_1 \mathcal{R}$ with

$$\mathcal{E}_0 = 2q\bar{q}\left( \bar{\epsilon} - \epsilon \right) + \epsilon$$
$$\mathcal{E}_1 = (q - \epsilon) - \left( \bar{q} - q \right)\left( \bar{\epsilon} - \epsilon \right)(q - \epsilon) - 2q\bar{q}\left( \bar{\epsilon} - \epsilon \right). \tag{64}$$

We also checked that the linear coefficient inside the parentheses in equation (63) is always numerically positive in the interval $\epsilon < q < 1/2$.

As the stall time $\tau_s$ is defined by $\eta(\tau_s) = 0$, equation (63) immediately gives us

$$e^{-\tau_s} \approx A_s \, \eta_{\max}^2, \tag{65}$$

**Figure 10.** (a) The efficiency $\eta$ and power gain $\langle P \rangle$ as functions of $e^{-\tau}$. (b) $e^{-\tau_s}$ and $e^{-\tau_{op}}$ as functions of $\eta_{max}$. In (b), the dotted line $\sim x^2$ is obtained from equation (65), and the dashed line corresponds to equation (71). We choose $t_R = t_M$, $q = 0.2$, $\epsilon = 0.1$, and $\beta_R = 1$ for both (a) and (b). We set $\beta_M = 5$ for (a) and vary $\beta_M$ for (b).



**Figure 11.** The efficiency $\eta_{op}$ at the optimal power for various $\epsilon$ and $q$. Here we choose $k = k'$ and $t_R = t_M$ for simplicity.

where the coefficient $A_s$ is given by

$$A_s = \left( \frac{\mathcal{E}_1}{\mathcal{E}_0 - \epsilon} + \frac{\mathcal{E}_0 - \epsilon}{q - \epsilon} \right)^{-2}. \tag{66}$$

It turns out that the scaling behavior in equation (65) extends quite well to finite $\eta_{\max}$, as can be seen in figure 10(b).

The optimal time $\tau_{op}$ should satisfy

$$\frac{\mathrm{d}\langle P\rangle}{\mathrm{d}\tau}\bigg|_{\tau_{op}} = 0,$$

which is rewritten as

$$\frac{\langle W_{\mathrm{net}}(\mathcal{R}_{op})\rangle}{\tau_{op}} = \frac{\mathrm{d}\langle W_{\mathrm{net}}(\mathcal{R})\rangle}{\mathrm{d}\tau}\bigg|_{\tau_{op}}, \tag{67}$$

where $\mathcal{R}_{op} = \mathrm{e}^{-\tau_{op}/2}$. As $\tau_{op} > \tau_s$, $\mathcal{R}_{op}$ should be also small. This allows us to expand the above equation for small $\mathcal{R}_{op}$, yielding

$$1 - \frac{\beta_R}{\beta_M}\lambda_{op} = \frac{\tau_{op}}{2}\left(\frac{\mathcal{E}_1}{\mathcal{E}_0 - \epsilon} + \frac{\mathcal{E}_0 - \epsilon}{q - \epsilon}\right)\mathcal{R}_{op}, \tag{68}$$

where $\mathcal{E}_0$ and $\mathcal{E}_1$ are the same as in equation (64), and $\lambda_{op}$ is given by

$$\lambda_{op} = \lambda_\infty + \lambda_\infty\left(\frac{\mathcal{E}_1}{\mathcal{E}_0 - \epsilon} + \frac{\mathcal{E}_0 - \epsilon}{q - \epsilon}\right)\mathcal{R}_{op}. \tag{69}$$

Plugging the above expression for $\lambda_{op}$ into equation (68), we get

$$A_s^{-1/2}\mathcal{R}_{op} = \frac{\eta_{\max}}{1 + \dfrac{\tau_{op}}{2}} = \frac{\eta_{\max}}{1 - \ln\mathcal{R}_{op}}, \tag{70}$$

which yields

$$\mathrm{e}^{-\tau_{op}} \approx A_s\left(\frac{\eta_{\max}}{\ln\eta_{\max}}\right)^2 \frac{1}{\left[1 + f\left(-\ln\eta_{\max}\right)\right]^2}, \tag{71}$$

where

$$f(x) = \frac{1}{x}\left[\ln x + 1 - \frac{1}{2}\ln A_s + \frac{\ln x}{x}\right]. \tag{72}$$

Finally, inserting equation (71) into equation (63), it is straightforward to find an approximate expression for $\eta_{op}$:

$$\eta_{op} \approx \left(1 - \frac{1}{|\ln\eta_{\max}|\left[1 + f\left(-\ln\eta_{\max}\right)\right]}\right)\eta_{\max}. \tag{73}$$

Therefore, in the limit of $\eta_{\max} \to 0$, one can see $\eta_{op} \approx \eta_{\max}$, not $\frac{1}{2}\eta_{\max}$, and that the next correction is logarithmic and therefore quite slow. This calculation confirms that the linear irreversible thermodynamics slightly out of equilibrium in [25] should *not* be applicable in our case, simply because our processes are far from equilibrium. Indeed, in our case, the entropy production is maximal in the limit of $\eta_{\max} \to 0$. In future studies, the validity of $\eta_{op} \approx \eta_{\max}$ should be addressed in the context of universality for general information engines, showing the maximum efficiency at the same point where the entropy production is maximal.

## 9. Conclusions

In this work, we have studied a simple example of an information engine which can be realized physically in terms of stochastic Markov processes. In agreement with previous studies, we find that the information feedback allows one to extract work in a situation where this would be thermodynamically impossible without feedback. Moreover, we confirm that total entropy production during relaxation obeys a fluctuation theorem, implying that the extracted work is bounded *from above* by the mutual information gain between memory and system.

Providing a physical realization of the memory and the feedback loop, we have shown that the measurement process (i.e., the information-processing part of the information engine) exhibits similar properties, which are opposite in character. In particular, the entropy production during measurement is found to obey a fluctuation theorem as well. This implies that the measurement process itself costs energy, and that this additional energy supply is bounded *from below* by the same mutual information gain. Putting these pieces together, it is no

surprise that the total setup consisting of system and memory satisfies the conventional second law of thermodynamics. Thus, we have shown that the thermodynamic second law, which is required to hold for the entire system during any finite process, leads to a duality in the properties of system and memory in this kind of information engine.

For simplicity, we have presented most of our analytic results in the limit of infinite measurement- and relaxation time. However, the extension to finite times is straightforward. At the end of the paper, we have explicitly described some numerical results for finite-time measurement and relaxation. As in conventional heat engines, the efficiency of the information engine is maximized when the cycle time becomes infinite. However, in contrast to conventional heat engines, the entropy production is also maximal in this limit. On the other hand, we have demonstrated that the power gain acquires its maximum at a finite cycle time. We have also discussed the relation between the maximal efficiency and the efficiency at the operating point of maximal power.

The striking differences between our model and conventional reversible heat engines can be traced back to the fact that our setup operates under non-equilibrium conditions. It would be interesting to investigate to what extent our observations can be explained in a universal framework.

## Acknowledgments

## References

[1] Earman J and Norton J D 1998 *Studies in the History and Philosophy of Modern Physics* **29** 435
     Earman J and Norton J D 1999 *Studies in the History and Philosophy of Modern Physics* **30** 1
[2] Debye P J W, Nernst W, Smoluchowski M, Sommerfeld A and Lorentz H A 1914 *Vorträge über die Kinetische Theorie der Materie und der Elektrizität* (Leipzig: Teubner) vol 6 p 89
[3] Szilard L 1929 *Z. Phys.* **53** 840
[4] Toyabe S, Sagawa T, Ueda M, Muneyuki E and Sano M 2010 *Nat. Phys.* **6** 988–92
[5] Landauer R 1961 *IBM J. Res. Dev.* **5** 183
[6] Bennett C H 1973 *IBM J. Res. Dev.* **17** 525
[7] Sagawa T and Ueda M 2008 *Phys. Rev. Lett.* **100** 080403
[8] Sagawa T and Ueda M 2010 *Phys. Rev. Lett.* **104** 090602
[9] Sagawa T and Ueda M 2012 *Phys. Rev. Lett.* **109** 180602
[10] Sagawa T and Ueda M 2012 *Phys. Rev.* E **85** 021104
[11] Horowitz J M, Sagawa T and Parrondo J M R 2013 *Phys. Rev. Lett.* **111** 010602
[12] Abreu D and Seifert U 2011 *Europhys. Lett.* **94** 10001
[13] Bauer M, Abreu D and Seifert U 2012 *J. Phys.* A **45** 162001
[14] Sandber H, Delvenne J-C, Newton N J and Mitter S K 2014 *Phys. Rev.* E **90** 042119
[15] Koski J V, Maisi V F, Sagwa T and Pekola J P 2014 *Phys. Rev. Lett.* **113** 030601
[16] Barato A C, Hartich D and Seifert U 2013 *J. Stat. Phys.* **153** 460
[17] Barato A C and Seifert U 2013 *Europhys. Lett.* **101** 60001
[18] Barato A C, Hartich D and Seifert U 2013 *Phys. Rev.* E **87** 042104
[19] Barato A C and Seifert U 2014 *Phys. Rev. Lett.* **112** 090601
[20] Strasberg P, Schaller G, Brandes T and Esposito M 2013 *Phys. Rev. Lett.* **110** 040601
[21] Schnakenberg J 1976 *Rev. Mod. Phys.* **48** 571
[22] Andrieux D and Gaspard P 2004 *J. Chem. Phys.* **121** 6167
[23] Seifert U 2005 *Phys. Rev. Lett.* **95** 040602
[24] Curzon F and Ahlborn B 1975 *Am. J. Phys.* **43** 22
[25] van den Broeck C 2005 *Phys. Rev. Lett.* **95** 190602
[26] Esposito M, Lindenberg K and van den Broeck C 2009 *Phys. Rev. Lett.* **102** 130602