



Contents lists available at ScienceDirect

# Journal of Computational and Applied Mathematics

journal homepage: [www.elsevier.com/locate/cam](http://www.elsevier.com/locate/cam)

## Representations and divergences in the space of probability measures and stochastic thermodynamics

Liu Hong<sup>a,b</sup>, Hong Qian<sup>a,\*</sup>, Lowell F. Thompson<sup>a</sup><sup>a</sup> Department of Applied Mathematics, University of Washington, Seattle, WA 98195-3925, USA<sup>b</sup> Zhou Pei-Yuan Center for Applied Mathematics, Tsinghua University, Beijing, 100084, PR China

### ARTICLE INFO

#### Article history:

Received 5 February 2019

Received in revised form 2 March 2020

#### MSC:

60-xx

80-xx

82-xx

#### Keywords:

Radon–Nikodym derivative

Affine structure

Space of probability measures

Heat divergence

### ABSTRACT

Radon–Nikodym (RN) derivative between two measures arises naturally in the affine structure of the space of probability measures with densities. Entropy, free energy, relative entropy, and entropy production as mathematical concepts associated with RN derivatives are introduced. We identify a simple equation that connects two measures with densities as a possible mathematical basis of the entropy balance equation that is central in nonequilibrium thermodynamics. Application of this formalism to Gibbsian canonical distribution yields many results in classical thermomechanics. An affine structure based on the canonical representation and two divergences are introduced in the space of probability measures. It is shown that thermodynamic work, as a conditional expectation, is indicative of the RN derivative between two energy representations being singular. The entropy divergence and the heat divergence yield respectively a Massieu–Planck potential based and a generalized Carnot inequalities.

© 2020 Elsevier B.V. All rights reserved.

## 1. Introduction

A subtle distinction exists between the prevalent approach to stochastic processes in traditional applied mathematics and the physicist's perspective on stochastic dynamics: In Kolmogorov's theory of stochastic processes, the dynamics are described in terms of a trajectory  $\{\mathbf{x}(t) : t \in [0, \infty)\}$ . Applied mathematicians treat each of these trajectories as a random event in a large probability space and then study probability distributions over the space of all possible trajectories  $\mathbf{x}(t)$ . Physicists, however, are more accustomed to thinking of a "probability distribution changing with time",  $\rho(\mathbf{x}, t)$ . In the case of continuous-time Markov processes,  $\rho(\mathbf{x}, t)$  is described by the solution to a Fokker–Planck equation or a master equation, while for a Markov chain simply by a stochastic matrix. This latter perspective can perhaps be more rigorously formulated in the *space of probability measures*. The dynamics are then represented as a *change of measure*. The Radon–Nikodym (RN) derivative is a key mathematical concept associated with changes in measures [1]. Interestingly, RN derivative between two measures is also at the heart of the concept of *fluctuating entropy* [2,3].

This "probability distribution changing with time" view is, of course, not foreign to mathematics. Actually in the 1950s, the stochastic diffusion process developed by Feller, Nelson, and others was precisely a such theory [4–7]. That approach, based on solutions to linear parabolic partial differential equations, was formulated in a linear function space. We now know that a more geometrically intrinsic representation for the space of probability measures cannot be linear: There is simply no natural choice of origin. Rather, an *affine space* is more appropriate [8,9].

\* Corresponding author.

E-mail addresses: [zcamhl@tsinghua.edu.cn](mailto:zcamhl@tsinghua.edu.cn) (L. Hong), [hqian@uw.edu](mailto:hqian@uw.edu) (H. Qian), [lthomps@uw.edu](mailto:lthomps@uw.edu) (L.F. Thompson).

Entropy and energy are key concepts in the classical theory of thermodynamics, which is now well understood to have a probabilistic basis. In fact, one could argue that the very notion of “heat” arises only when one treats the motions of deterministic Newtonian point masses as stochastic. In the statistical treatment of thermodynamics, Gibbs’ canonical energy distribution is one of the key results that characterize a thermodynamic equilibrium [10]. As we shall see, it figures prominently in the affine space.

The foregoing discussion suggests the possibility of re-thinking thermodynamics and information theory in a novel mathematical framework [11]. Both information theory and thermodynamics are concerned with notions such as entropy, free energy and relative entropy. These concepts are introduced in Section 2 under a single framework based on the Radon–Nikodym derivative, as a random variable relating two different measures. In its broadest context, we are able to capture the essential mathematics used in the theory of equilibrium and nonequilibrium thermodynamics. This approach significantly enriches the scope of “information theory” [12]. The RN derivative should not be treated as an esoteric mathematical concept: It is simply a powerful way to quantify even infinitesimal changes in the probability distributions; it is the calculus for thinking of change in terms of chance [13].

In Section 3, the notion of a temperature,  $T = \beta^{-1}$  is introduced through the canonical probability distribution  $Z^{-1}(\beta)e^{-\beta U(\omega)}$ . It has been shown recently that this Gibbsian distribution has a much broader applications than just thermal physics: It is in fact a limit theorem of a sequence of conditional probability densities under an additive quasi-conservative observable [14]. The focus of this section is to show the centrality of RN derivative in the theory of thermodynamics. The RN derivative is used to describe several results in physics that includes the thermodynamic cycle, equation of states, and the Jarzynski–Crooks equalities.

Next, in Section 4 we equip the space of probability measures with an affine structure and show that the canonical distribution with a random variable  $U(\omega)$  and a parameter  $\beta$  becomes precisely an affine line in the space of probability measures when one particular measure  $\mathbb{P}$  is chosen as a reference point. With this, the tangent space becomes a linear vector space of random variables and it provides a representation for the space of probability measures. A series of results are obtained. Readers who are more mathematically inclined can skip Section 3, come directly to Section 4, and then go back to Section 3 afterward.

Section 5 contains some discussions.

The presentation of the paper is not mathematically rigorous. The emphasis is on illustrating how the pure mathematical concepts can be fittingly applied in narrating this branch of physics. More thorough treatments of the subject are forthcoming [9].

## 2. Entropy, relative entropy, and a fundamental equation of information

### 2.1. Information and entropy

Information theory owes (to a large extent) its existence as a separate subject from the theories of probability and statistics to a singular emphasis on the notion of *entropy* as a quantitative measure of information. It is important to point out at the outset that *information* is a random variable, defined on a probability space  $(\Omega, \mathcal{F}, \mathbb{P})$ , through a Radon–Nikodym derivative  $\frac{d\mathbb{P}}{d\mu}(\omega)$ ,  $\omega \in \Omega$ , between two measures  $\mathbb{P}$  and  $\mu$  that are absolutely continuous w.r.t. each other [2,3,15]. If the  $\Omega \subseteq \mathbb{R}^n$  and  $\mu$  is the Lebesgue measure, then

$$-\ln \left( \frac{d\mathbb{P}}{d\mu}(\omega) \right) \tag{1}$$

is the *self-information* [16,17], which is a random variable and its expected value is the standard form of Shannon entropy:

$$S[\mathbb{P}] \triangleq - \int_{\Omega} f(x) \ln f(x) dx, \tag{2}$$

in which the Radon–Nikodym derivative is the probability density function,  $\frac{d\mathbb{P}}{d\mu} \equiv f(x)$ .

In general, if  $\mu$  is normalizable, then one has a maximum entropy inequality  $S[\mathbb{P}] \leq \ln \mu(\Omega) < +\infty$ . Similarly, one has the free energy

$$H[\mathbb{P} \parallel \mu] \triangleq \int_{\Omega} \ln \left( \frac{d\mathbb{P}}{d\mu}(\omega) \right) d\mathbb{P}(\omega) \geq -\ln \mu(\Omega). \tag{3}$$

When  $\mu$  is also a normalized probability measure  $\mathbb{P}'$ , the  $H[\mathbb{P} \parallel \mathbb{P}']$  is called the *relative entropy* or Kullback–Leibler (KL) divergence. The minimum free energy inequality in (3) becomes the better known, but less interesting,  $H[\mathbb{P} \parallel \mathbb{P}'] \geq 0$ .

From now on, we will drop most references to the underlying space  $(\Omega, \mathcal{F})$ . Moreover, we will assume that  $\Omega \subseteq \mathbb{R}^n$  with the usual  $\sigma$ -algebra and that  $\mathbb{P}$  is absolutely continuous w.r.t. the Lebesgue measure. These conditions are not strictly necessary, but they simplify the notation considerably in illustrating our key ideas.

### 2.2. Fundamental equation of information

With the various forms of entropy introduced above and some straightforward statistical logic, one naturally has the following equation that involves three measures: two probabilistic and the Lebesgue. In particular, let  $\mathbb{P}_1$  and  $\mathbb{P}_2$  be two probability measures with density functions  $f_1(x)$  and  $f_2(x)$  with respect to the Lebesgue measure:

$$\begin{aligned} \Delta S &= S[\mathbb{P}_2] - S[\mathbb{P}_1] = \int_{\mathbb{R}} f_1(x) \ln f_1(x) dx - \int_{\mathbb{R}} f_2(x) \ln f_2(x) dx \\ &= \underbrace{\int_{\mathbb{R}} f_1(x) \ln \left( \frac{f_1(x)}{f_2(x)} \right) dx}_{\Delta S^{(i)}: \text{entropy production}} + \underbrace{\int_{\mathbb{R}} (f_2(x) - f_1(x)) (-\ln f_2(x)) dx}_{\Delta S^{(e)}: \text{entropy exchange}}. \end{aligned} \tag{4}$$

The entropy production  $\Delta S^{(i)}$  is never negative, while the entropy exchange  $\Delta S^{(e)}$  has no definitive sign. If  $f_2(x)$  is the unique invariant density of some measure-preserving dynamics [18], then  $-\ln f_2(x)$  is customarily referred to as the “equilibrium energy function”, then  $\Delta S^{(e)}$  is the change in the “mean energy”, which is related to “heat”.

Entropy and free energy in (2) and (3) have their namesakes in the theory of statistical equilibrium thermodynamics [10]. The Second Law, in terms of entropy maximization or free energy minimization, has its statistical basis precisely in the two inequalities associated with  $S$  and  $H$ . The  $\Delta S^{(i)}$  term on the rhs of (4), however, is a nonequilibrium free energy associated with a *nonequilibrium distribution*, either due to a spontaneous fluctuation or a man-made perturbation [19]. In the theory of stochastic dynamics, one uses a probability distribution  $\rho(x, t)$  to represent the state of a system; thus any  $\rho$  that differs from the equilibrium distribution is a nonequilibrium distribution. In applications to laboratory systems, the  $\rho$  can only be obtained from a data-based statistical approach. This approach can rely on either a time scale separation, or a system of many independent and identically distributed subsystems, or a fictitious ensemble. Ideal gas theory and the Rouse model of polymers are two successful examples of the second type [19].

Eq. (4) in fact has the form of the *fundamental equation of nonequilibrium thermodynamics*. It states that if  $f_2(x)$  is uniform, then  $\Delta S = \Delta S^{(i)} \geq 0$ ; and if one identifies  $U(x) \triangleq -T \ln f_2(x)$ , where  $T$  is a positive constant, then one can introduce  $F[\mathbb{P}] \triangleq \mathbb{E}^{\mathbb{P}}[U] - TS[\mathbb{P}]$ , and  $\Delta F = T \Delta S^{(i)} \geq 0$ . Unifying the various forms of the Second Law to a single concept of entropy production was a key idea of the Brussel school of thermodynamics [20].<sup>1</sup> See [3,11,21,22], and the references cited within, for the theory of entropy production of Markov processes.

### 2.3. Two results on relative entropy

With regards to relative entropy, there are two results worth discussing.

First, as the expected value of the logarithm of the Radon–Nikodym derivative  $\xi \equiv \ln \left( \frac{d\mathbb{P}_1}{d\mathbb{P}_2}(\omega) \right)$ , the relative entropy between two probability measures can be written as

$$H[\mathbb{P}_1 \parallel \mathbb{P}_2] = \int_{\mathbb{R}} f_1(x) \ln \left( \frac{f_1(x)}{f_2(x)} \right) dx = \mathbb{E}^{\mathbb{P}_1}[\xi(\omega)], \tag{5}$$

with respective probability density functions  $f_1(x) = \frac{d\mathbb{P}_1(x)}{dx}$  and  $f_2(x) = \frac{d\mathbb{P}_2(x)}{dx}$ . The non-negativity of the  $H[\mathbb{P}_1 \parallel \mathbb{P}_2]$  can actually be framed as a consequence of a stronger result, an equality

$$\mathbb{E}^{\mathbb{P}_1} [e^{-\xi(\omega)}] = 1, \tag{6}$$

and an inequality for convex exponential function:

$$\mathbb{E}^{\mathbb{P}_1} [\xi(\omega)] \geq -\ln \mathbb{E}^{\mathbb{P}_1} [e^{-\xi(\omega)}] = 0. \tag{7}$$

Eq. (6) implies that the Second Law and entropy production could even be formulated through equalities rather than inequalities. Indeed, variations of (6) have found numerous applications in thermodynamics, such as Zwanzig’s free energy perturbation method [23], the Jarzynski–Crooks relation [24,25], and the Hatano–Sasa equality [26].

Second, if the density  $f_2$  contains an unknown parameter  $\theta$ , then  $f_2(x; \theta)$  is the likelihood function for  $\theta$ . In this case, with respect to the change of measure,

$$\begin{aligned} \mathcal{I}_\ell(\theta) &\triangleq -\mathbb{E}^{\mathbb{P}_2} \left[ \frac{\partial^\ell}{\partial \theta^\ell} \ln f_2(\omega; \theta) \Big| \theta \right] \\ &= -\int_{\mathbb{R}} f_2(x; \theta) \frac{\partial^\ell}{\partial \theta^\ell} \ln f_2(x; \theta) dx \end{aligned}$$

<sup>1</sup> The second author would like to acknowledge an enlightening discussion with M. Esposito in the spring of 2011 at the Snogeholm Workshop on Thermodynamics, Sweden.

$$\begin{aligned}
 &= - \int_{\mathbb{R}} \left\{ \left( \frac{f_2(x; \theta)}{f_1(x)} \right) \frac{\partial^\ell}{\partial \theta^\ell} \ln \left( \frac{f_2(x; \theta)}{f_1(x)} \right) \right\} f_1(x) dx \\
 &= \mathbb{E}^{\mathbb{P}^1} \left[ e^{-\xi(\omega)} \frac{\partial^\ell}{\partial \theta^\ell} \xi(\omega; \theta) \right].
 \end{aligned} \tag{8}$$

$\mathcal{I}_0(\theta)$  is the Shannon entropy of  $X_2(\theta)$ ,  $\mathcal{I}_1(\theta) \equiv 0$ , and  $\mathcal{I}_2(\theta)$  is the Fisher Information for  $X_2(\theta)$ :

$$\mathcal{I}_2(\theta) = \mathbb{E} \left[ \left( \frac{\partial}{\partial \theta} \ln f_2(X_2; \theta) \right)^2 \middle| \theta \right]. \tag{9}$$

### 3. Canonical distribution and thermodynamics

In many applications, stochastic dynamics exhibit a separation of slow and fast time scales [27,28]. In mechanical systems with sufficiently small friction, the dynamics are organized as fast Hamiltonian dynamics with slow energy dissipation through heat. The theory of thermodynamics arises in this context when the mechanical motions of point masses are described stochastically. It can be shown then that the probability distribution for the energy  $E$  of a small mechanical system in equilibrium with a large heat bath takes a particularly canonical form

$$p_E(y) = \frac{\Omega^{(B)}(y)e^{-\beta y}}{Z(\beta)}, \tag{10}$$

in which  $\beta^{-1} = T$  is the temperature of the heat bath [10,14]. In fact, if  $\mathbf{x}$  denotes random variable in an appropriate state space and  $U(x)$  is the mechanical energy function, then one has distribution  $f_{\mathbf{x}}(x) \propto e^{-\beta U(x)}$ , and

$$p_E(y) dy = \int_{y < U(x) \leq y+dy} f_{\mathbf{x}}(x) dx = \left( \frac{\Omega^{(B)}(y)e^{-\beta y}}{Z(\beta)} \right) dy, \tag{11}$$

in which

$$\Omega^{(B)}(y) = \frac{1}{dy} \int_{y < U(x) \leq y+dy} dx = \frac{d\Omega^{(G)}(y)}{dy}, \quad \Omega^{(G)}(y) = \int_{U(x) \leq y} dx. \tag{12}$$

In  $\Omega^{(B)}$  and  $\ln \Omega^{(G)}$  are called Boltzmann’s entropy and Gibbs’ entropy in statistical mechanics [29]. They are related via  $d\Omega^{(G)}(y) = \Omega^{(B)}(y)dy$ . That is,  $\Omega^{(G)}$  is a cumulative distribution function and  $\Omega^{(B)}$  is its density function.

Note that the expected value of any function of the energy  $U(x)$  (e.g.,  $g(U)$ ) is invariant under different representations as a result of the rules of changes of variable for integration. For example, if  $\mathbf{x}$  is a state space representation and  $E$  is the energy representation, then

$$\begin{aligned}
 \int_{\mathbb{R}} g(U(x))f_{\mathbf{x}}(x)dx &= \int_{\mathbb{R}} g(y) \left( \frac{e^{-\beta y}}{Z(\beta)} \right) \Omega^{(B)}(y)dy \\
 &= \int_{\mathbb{R}} g(y)p_E(y)dy.
 \end{aligned}$$

In contrast, the thermodynamic entropy in statistical mechanics is not invariant under different representations [30]:

$$- \int p_{\mathbf{x}}(x) \ln p_{\mathbf{x}}(x) dx = - \int_{\mathbb{R}} p_E(y) \ln \left( \frac{p_E(y)}{\Omega^{(B)}(y)} \right) dy \tag{13a}$$

$$\neq - \int_{\mathbb{R}} p_E(y) \ln p_E(y) dy. \tag{13b}$$

The rhs of (13a) is precisely the negative free energy with non-normalized  $\Omega^{(G)}(y)$  as the reference measure (which has density  $\Omega^{(B)}(y)$ ). The missing term from (13a) to (13b)

$$\int_{\mathbb{R}} p_E(y) \left( - \ln \Omega^{(B)}(y) \right) dy. \tag{14}$$

is contributed by the reference measure. It is mean-internal-energy like. We see that while  $\ln \Omega^{(B)}(y)$  is widely considered as an “entropic” term, it actually plays the role of an energetic term in the energy representation in (13a). In terms of this measure-theoretic framework, the distinction between entropy and energy is always relative. This has long been understood in the work of J. G. Kirkwood on the *potential of mean force*, which is itself temperature dependent [31].

#### 3.1. Thermodynamics under a single temperature

**Equilibrium statistical thermodynamics.** In terms of the canonical distribution in (10), an equilibrium system under a constant temperature  $T = \beta^{-1}$  has its mechanical energy distributed according to the canonical distribution  $p^{eq}(y) =$

$Z^{-1}(\beta)\Omega(y)e^{-\beta y}$ . (We have dropped the superscript in  $\Omega^{(B)}(y)$  to avoid cluttering.) The mean internal energy associated with the  $p^{eq}(y)$  is then the expected value

$$\bar{U}(\beta) = \int_{\mathbb{R}} y \left( \frac{\Omega(y)e^{-\beta y}}{Z(\beta)} \right) dy = - \frac{d \ln Z(\beta)}{d\beta}, \tag{15a}$$

which can be decomposed into an equilibrium free energy and an entropy,  $\bar{U}(\beta) = F^{eq}(\beta) + \beta^{-1}S(\beta)$ , where:

$$\underbrace{F^{eq}(\beta) = -\beta^{-1} \ln Z(\beta)}_{\text{free energy}} \text{ and } \underbrace{S(\beta) = - \left( \frac{dF^{eq}(\beta)}{d\beta^{-1}} \right)}_{\text{entropy}}. \tag{15b}$$

One can verify that the  $S(\beta)$  is the same as (13a), but not (13b).

**Nonequilibrium statistical thermodynamics.** For deep mathematical reasons that will become clear in Section 4, discussions of nonequilibrium systems should begin in the full state space. Intuitively, the canonical energy representation  $p^{eq}(E)$  based on a given energy function  $U(x)$  is a “projection” in the space of probability measures that is nonholographic.

Consider a system outside statistical equilibrium with a nonequilibrium probability measure  $\mu^{neq}$ . Suppose that this measure is absolutely continuous w.r.t. some other probability measure  $\mathbb{P}$ , with density  $\rho(x) = \frac{d\mu^{neq}}{d\mathbb{P}}(x)$ . The measure  $\mu^{neq}$  possesses a nonequilibrium free energy functional (a potential that can cause change) given by

$$F^{neq}[\rho; \beta] \triangleq F^{eq}(\beta) + \beta^{-1} \int_{\Omega} \rho(x) \ln \left( \frac{\rho(x)}{p^{eq}(x)} \right) dx \tag{16a}$$

$$= \beta^{-1} \int_{\Omega} \rho(x) \ln \left( \frac{\rho(x)}{e^{-\beta U(x)}} \right) dx. \tag{16b}$$

One should recognize the fraction in (16b) as a Radon–Nikodym derivative of  $\rho$  w.r.t. the non-normalized canonical equilibrium measure  $e^{-\beta U(x)}$ . The minimum free energy inequality in (3) takes the form  $F^{neq}[\rho; \beta] \geq F^{eq}(\beta)$  for any distribution  $\rho$ . In fact,  $\beta\{F^{neq}[\rho; \beta] - F^{eq}(\beta)\}$  is the entropy production associated with the spontaneous relaxation process of the distribution  $\rho$  tending to  $p^{eq}$ .

The  $F^{neq}[\rho; \beta]$  also has another expression:

$$F^{neq}[\rho; \beta] = \underbrace{\beta^{-1} \int_{\Omega} \rho(x) \ln \rho(x) dx}_{\text{neg-entropy}} + \int_{\Omega} \rho(x) \underbrace{\left( -\beta^{-1} \ln p^{eq}(x) + F^{eq}(\beta) \right)}_{\text{internal energy of state } x, F^{eq} \text{ as reference}} dx. \tag{17}$$

Eq. (17) is very telling: The internal energy of a system in state  $x$  is given in the second term with a fixed energy gauge (i.e., the arbitrary constant in the  $U(x)$ ) according to the equilibrium  $F^{eq}$ , where  $U(x) = F^{eq}(\beta) - \beta^{-1} \ln p^{eq}(x)$ . This fact implies that a change in the energy function from  $U_1(x)$  to  $U_2(x)$  necessarily involves a change of gauge. Mechanical work in classical thermodynamics can be understood as a consequence of gauge invariance. One particular  $\beta$  defines an autonomous, time-homogeneous stochastic dynamical system with a unique  $p^{eq}$ . All the energetic discussions in such a system are with respect to the equilibrium free energy  $F^{eq}(\beta)$ , which fixes a choice for the energy gauge. In the theory of probability, the gauge invariance is achieved through the notion of conditional probability and the law of total probability.

### 3.2. A clarification of Eq. (16)

A discussion on the meaning of the expression in (16) is in order. To do that, let us only consider discrete  $x_k$ , and the corresponding

$$F^{neq}[\rho; \beta] = \sum_k \rho(x_k) \left[ \beta^{-1} \ln \left( \frac{\rho(x_k)}{p^{eq}(x_k)e^{-\beta F^{eq}(\beta)}} \right) \right]. \tag{18}$$

For a particular state  $z$ , if  $\rho(x) = \delta_{x,z}$ , then  $F^{neq}[\rho; \beta] = F^{eq}(\beta) - \beta^{-1} \ln p^{eq}(z)$ , which represents the traditional potential energy of the system in the state  $z$ . A question then naturally arises: Why is  $F^{neq}[\rho; \beta]$  the average of

$$\beta^{-1} \ln \left( \frac{\rho(x_k)}{p^{eq}(x_k)e^{-\beta F^{eq}}} \right), \tag{19a}$$

but not

$$\beta^{-1} \ln \left( \frac{1}{p^{eq}(x_k)e^{-\beta F^{eq}}} \right)? \tag{19b}$$

Actually, (19b) is the potential energy for a deterministic initial state  $x_k$ . It is natural, therefore, the average would be carried out over (19b) if the initial state of the system were a mixture of heterogeneous states (mhs). However, if the initial state is a stochastic fluctuating state (sfs), then the entropy of assimilation applies [32] and the  $F^{neq}[\rho; \beta]$  in (16a) is the

average carried out over (19a). The change from mhs to sfs is analogous to a change from the Lagrangian to the Eulerian representation in fluid mechanics; in stochastic terms, the potential for an sfs to do work is lower than an mhs [33].

### 3.3. Work, heat, and Jarzynski–Crooks’ relation

We now consider the case where the distribution  $\rho(x)$  in (16) arises from the equilibrium distribution  $p^{eq}(x)$  as the consequence of a temperature change from  $T_a$  to  $T_b$ :  $\rho(x) = Z^{-1}(\beta_a)e^{-\beta_a U(x)}$ , and the  $p^{eq}(x) = Z^{-1}(\beta_b)e^{-\beta_b U(x)}$ . Note that in the energy representation they can be written as  $\rho_E(y) = Z^{-1}(\beta_a)\Omega(y)e^{-\beta_a y}$  and  $p_E^{eq}(y) = Z^{-1}(\beta_b)\Omega(y)e^{-\beta_b y}$ ; they share the same Gibbs entropy  $\ln \Omega(y)$  determined by  $U(x)$  as in (12). Then

$$F^{neq}[\rho; \beta_b] - F^{eq}(\beta_b) = \beta_b^{-1} \int_{\Omega} \rho(x) \ln \left( \frac{\rho(x)}{p^{eq}(x)} \right) dx \tag{20a}$$

$$= \beta_b^{-1} \int_{\mathbb{R}} \rho_E(y) \ln \left( \frac{\rho_E(y)}{p_E^{eq}(y)} \right) dy \tag{20b}$$

$$= [\bar{U}(\beta_a) - \beta_b^{-1}S(\beta_a)] - [\bar{U}(\beta_b) - \beta_b^{-1}S(\beta_b)]. \tag{20c}$$

The equation from (20a) to (20b) utilizes a key property of a Radon–Nikodym derivative: *When it exists, it is invariant under a change of measure.*

Eq. (20c) is not widely discussed, but it is a highly meaningful result. It contains the essence of Crooks’ equality in time-inhomogeneous Markov processes [25]. It implies that at the instant of switching from  $T_a$  to  $T_b$ , the system has internal energy  $\bar{U}(\beta_a)$ , entropy  $S(\beta_a)$ , and nonequilibrium free energy

$$F^{neq}[\rho; \beta] = \bar{U}(\beta_a) - T_b S(\beta_a). \tag{21}$$

Assuming that both  $\rho(x)$  and  $p^{eq}(x)$  have the same  $\Omega(y)$ , Eq. (20) gives the free energy change that is expected to be the maximum reversible work that can be extracted. We now explicitly consider a change from  $\rho(x)$  to  $p^{eq}(x)$  that involves changing the mechanical energy function from  $U_1(x)$  to  $U_2(x)$ . Even though the corresponding canonical energy distributions are  $\rho_E(y) = Z_1^{-1}(\beta_a)\Omega_1(y)e^{-\beta_a y}$  and  $p_E^{eq}(y) = Z_2^{-1}(\beta_b)\Omega_2(y)e^{-\beta_b y}$ , these RN derivative  $\frac{d\rho_E}{dp^{eq}}(\omega)$  can be infinity! Thus in this case one has to start with the full distributions on the state space:

$$\begin{aligned} \beta_b^{-1} \int_{\Omega} \rho(x) \ln \left( \frac{\rho(x)}{p^{eq}(x)} \right) dx &= [\bar{U}_1(\beta_a) - \bar{U}_2(\beta_b)] - \beta_b^{-1} [S_1(\beta_a) - S_2(\beta_b)] \\ &+ \int_{\Omega} \rho(x) [U_2(x) - U_1(x)] dx. \end{aligned} \tag{22}$$

The last term in (22) is identified as the irreversible work associated with the isothermal relaxation process with mechanical change from  $U_1(x)$  to  $U_2(x)$ ,

$$\bar{\mathcal{W}}_{12}(\beta_a) = \int_{\Omega} \rho(x) \mathcal{W}_{12}(x) dx, \tag{23}$$

in which  $\mathcal{W}_{12}(x)$  should be considered as the logarithm of the Radon–Nikodym derivative between two non-normalized measures

$$\mathcal{W}_{12}(x) = \beta_a^{-1} \ln \left( \frac{e^{-\beta_a U_1(x)}}{e^{-\beta_a U_2(x)}} \right) = \beta_b^{-1} \ln \left( \frac{e^{-\beta_b U_1(x)}}{e^{-\beta_b U_2(x)}} \right). \tag{24}$$

$\mathcal{W}_{12}(x)$  is actually not a function of  $\beta$ ; work done in an isothermal process is independent of the temperature. In the canonical energy representation of  $U_1(x)$ , then,

$$\begin{aligned} \bar{\mathcal{W}}_{12}(\beta_a) &= \int_{\Omega} \rho(x) [U_2(x) - U_1(x)] dx \\ &= \int_{\mathbb{R}} \left( \frac{\Omega_1(y)e^{-\beta_a y}}{Z_1(\beta_a)} \right) \left\{ \frac{\int_{y < U_1(x) \leq y+dh} U_2(x) dx}{\int_{y < U_1(x) \leq y+dh} dx} - y \right\} dy. \end{aligned} \tag{25}$$

The first term inside  $\{ \dots \}$  is a conditional expectation:  $\mathbb{E}^{eq}[U_2(x)|U_1(x) = y]$ , where  $\mathbb{E}^{eq}$  is the expectation in terms of the equilibrium measure  $p^{eq}(x)$ .

The transferred irreversible heat is

$$\mathcal{Q}(\beta_b) \triangleq \beta_b^{-1} \left\{ S_1(\beta_a) - S_2(\beta_b) + \int_{\Omega} \rho(x) \ln \left( \frac{\rho(x)}{p^{eq}(x)} \right) dx \right\}. \tag{26}$$

Then the relation

$$S_2(\beta_b) - S_1(\beta_a) + \frac{Q(\beta_b)}{T_b} = \Delta S^{(i)} = \int_{\Omega} \rho(x) \ln \left( \frac{\rho(x)}{p^{eq}(x)} \right) dx \geq 0 \tag{27}$$

is known as the Clausius inequality in thermodynamics. The equality is a special case of the fundamental equation of nonequilibrium thermodynamics.

Concerning the work  $\mathcal{W}_{12}(x)$  in (24), we have Jarzynski–Crooks’ relation [24,25]:

$$\int_{\Omega} \left( \frac{e^{-\beta_a U_1(x)}}{Z_1(\beta_a)} \right) e^{-\beta_a \mathcal{W}_{12}(x)} dx = \int_{\Omega} \frac{e^{-\beta_a U_2(x)}}{Z_1(\beta_a)} dx = \frac{Z_2(\beta_a)}{Z_1(\beta_a)}. \tag{28}$$

Note that the work is performed under  $\beta_b$ , but the rhs of (28) is evaluated at  $\beta_a$ . The original Jarzynski–Crooks’ equality emphasized path-wise average over a stochastic trajectory, but Eq. (28) is an ensemble average over a single step, which can be generalized to many different other forms [34].

**The concept of exergy.** In Eq. (21), equilibrium internal energy and entropy under temperature  $T_a$ ,  $\bar{U}(T_a)$  and  $S(T_a)$  are assembled with temperature  $T_b \neq T_a$  to form a nonequilibrium free energy  $F^{neq} = \bar{U}(T_a) - T_b S(T_a)$ , which plays a central role in our analysis of canonical systems. This quantity has been extensively discussed in the literature on thermodynamics: *Exergy* of a system is “the maximum fraction of an energy form which can be transformed into work”. The remaining part is the waste heat [35]. After a system reaches equilibrium with its surrounding, its exergy is zero. Therefore, the concept of exergy epitomizes a nonequilibrium quantity [36]. Its identification to the entropy production in Eq. (20) implies its importance in information energetics. Even though the term “exergy” was coined as late as in 1956, the idea had been already in the work of Gibbs.

**Mechanical work of an ideal gas.** For an ideal gas with total mechanical energy  $U(x) = U_p(x_1) + U_k(x_2)$ , where  $U_p$  and  $U_k$  are potential and kinetic energy functions, and  $\mathbf{x}_1$  and  $\mathbf{x}_2$  are position and momentum state variables,

$$U(x) = \sum_{i=1}^N \left\{ \frac{x_{2,i}^2}{2m_i} + H_V(x_{1,i}) \right\}, \tag{29}$$

in which  $H_V(z) = 0$  when  $0 < z < V$  and  $H_V(z) = +\infty$  when  $z \leq 0$  or  $z \geq V$ . The  $V$  represents the “volume” of a box containing the ideal gas. Then

$$\Omega(E, V) = \frac{V^N}{dE} \int_{E < U_k(x_2) \leq E+dE} dx_2 = V^N \tilde{\Omega}(E, N), \tag{30}$$

in which the  $\tilde{\Omega}$  is independent of  $V$ . Therefore, the mechanical work associated with a change in  $V_1 = V \rightarrow V_2 = V + \Delta V$  is given by

$$\beta^{-1} \ln \left( \frac{\Omega(E, V_2)}{\Omega(E, V_1)} \right) = NT \ln \left( \frac{V + \Delta V}{V} \right) = \frac{NT \Delta V}{V} = \hat{p} \Delta V, \tag{31}$$

where  $\hat{p} = Nk_B T/V$  is the pressure of an ideal gas. (We have set Boltzmann’s constant  $k_B \equiv 1$  throughout the present paper.)

### 3.4. Application to heat engines and thermodynamic cycles

**Carnot cycle.** Applying Eqs. (24) and (26) twice for *thermomechanical* (i.e., temperature and mechanical) changes from  $\{T_a, U_1\}$  to  $\{T_b, U_2\}$  and from  $\{T_b, U_2\}$  back to  $\{T_a, U_1\}$ , we derive the celebrated Carnot efficiency for a heat engine. For each of the processes described in the left column below, the energetic status of the system is shown in the right column:

$$\text{adiabatic switching } \{T_a, U_1\} \rightarrow \{T_b, U_1\}: F_1^{neq}(T_b) = \bar{U}_1(T_a) - T_b S_1(T_a), \tag{32a}$$

$$\text{isothermal relaxation } \{T_b, U_1\} \rightarrow \{T_b, U_2\}: \bar{U}_1(T_a) - \bar{U}_2(T_b) = Q_{12}(T_b) - \bar{W}_{12}, \tag{32b}$$

$$\text{equilibrium under } T_b: F_2^{eq}(T_b) = \bar{U}_2(T_b) - T_b S_2(T_b), \tag{32c}$$

$$\text{adiabatic switching } \{T_b, U_2\} \rightarrow \{T_a, U_2\}: F_2^{neq}(T_a) = \bar{U}_2(T_b) - T_a S_2(T_b), \tag{32d}$$

$$\text{isothermal relaxation } \{T_a, U_2\} \rightarrow \{T_a, U_1\}: \bar{U}_2(T_b) - \bar{U}_1(T_a) = Q_{21}(T_a) - \bar{W}_{21}, \tag{32e}$$

$$\text{equilibrium under } T_a: F_1^{eq}(T_a) = \bar{U}_1(T_a) - T_a S_1(T_a). \tag{32f}$$

In (32f), the system is returned to the equilibrium state under  $T_a$ . Without loss of generality, let  $T_a > T_b$ . In the ideal Carnot cycle, one assumes that the processes of switching the temperatures are adiabatic without free energy dissipation. That is, the  $F_1^{neq}(T_b)$  in (32a) is strictly equal to  $F_1^{eq}(T_a)$  in (32f), with a reversible change of gauge reference, and similarly

the  $F_2^{\text{neq}}(T_a)$  in (32d) is strictly equal to  $F_2^{\text{eq}}(T_b)$  in (32c). In the two processes of isothermal relaxation, irreversible heat  $Q_{12}(T_b) = T_b\{S_1(T_a) - S_2(T_b) + \Delta S_{12}^{(i)}\}$  and  $Q_{21}(T_a) = T_a\{S_2(T_b) - S_1(T_a) + \Delta S_{21}^{(i)}\}$  each contain an entropy production term,

$$\Delta S_{jk}^{(i)} = \int_{\mathbb{R}} p_j^{\text{eq}}(y) \ln \left( \frac{p_j^{\text{eq}}(y)}{p_k^{\text{eq}}(y)} \right) dy \geq 0. \tag{33}$$

In a Carnot cycle with quasi-static processes, they are assumed to be zero. Then, the total work done by the system over the cycle is

$$\begin{aligned} W &= -(\overline{W}_{12} + \overline{W}_{21}) = -Q_{12}(T_b) - Q_{21}(T_a) \\ &= T_b \left\{ S_2 - S_1 - \int_{\mathbb{R}} p_1^{\text{eq}} \ln \left( \frac{p_1^{\text{eq}}}{p_2^{\text{eq}}} \right) dy \right\} + T_a \left\{ S_1 - S_2 - \int_{\mathbb{R}} p_2^{\text{eq}} \ln \left( \frac{p_2^{\text{eq}}}{p_1^{\text{eq}}} \right) dy \right\} \\ &\leq T_b [S_2(T_b) - S_1(T_a)] + T_a [S_1(T_a) - S_2(T_b)], \end{aligned} \tag{34}$$

in which the reversible heat being absorbed at  $T_a$  is  $Q_h = T_a[S_1(T_a) - S_2(T_b)] > 0$ , and the heat being expelled at  $T_b$  is  $Q_l = T_b[S_2(T_b) - S_1(T_a)] < 0$ . Thus the Carnot (first-law) efficiency

$$\eta_{\text{Carnot}} = \frac{W}{Q_h} \leq 1 - \frac{T_b}{T_a}. \tag{35}$$

On the other hand, since the rhs of (34) is the maximum possible work, the second-law, exergy efficiency

$$\eta_{\text{exergy}} = \frac{W}{(T_a - T_b)[S_1(T_a) - S_2(T_b)]} = \frac{W}{Q_h \left( 1 - \frac{T_b}{T_a} \right)} \leq 1. \tag{36}$$

**Stirling cycle.** There are many different realizations of heat engines in terms of thermodynamic cycles. We now consider the Stirling cycle below.

$$\text{isothermal working } \{T_a, U_1\} \rightarrow \{T_a, U_2\}: \quad \overline{U}_1(T_a) - \overline{U}_2(T_a) = \overline{Q}_{12}(T_a) - \overline{W}_{12}, \tag{37a}$$

$$\text{isochoric cooling } \{T_a, U_2\} \rightarrow \{T_b, U_2\}: \quad \overline{U}_2(T_a) - \overline{U}_2(T_b) = Q_2(T_a, T_b), \tag{37b}$$

$$\text{equilibrium under } \{T_b, U_2\}: \quad F_2^{\text{eq}}(T_b) = \overline{U}_2(T_b) - T_b S_2(T_b), \tag{37c}$$

$$\text{isothermal working } \{T_b, U_2\} \rightarrow \{T_b, U_1\}: \quad \overline{U}_2(T_b) - \overline{U}_1(T_b) = \overline{Q}_{21}(T_b) - \overline{W}_{21}, \tag{37d}$$

$$\text{isochoric heating } \{T_b, U_1\} \rightarrow \{T_a, U_1\}: \quad \overline{U}_1(T_b) - \overline{U}_1(T_a) = Q_1(T_b, T_a), \tag{37e}$$

$$\text{equilibrium under } \{T_a, U_1\}: \quad F_1^{\text{eq}}(T_a) = \overline{U}_1(T_a) - T_a S_1(T_a). \tag{37f}$$

After two isothermal processes in (37a), (37d), the system is still in the equilibrium states with free energy  $F_2^{\text{eq}}(T_a) = \overline{U}_2(T_a) - T_a S_2(T_a)$  and  $F_1^{\text{eq}}(T_b) = \overline{U}_1(T_b) - T_b S_1(T_b)$  respectively. Notice the difference between the equilibrium free energy above and the non-equilibrium free energy functions  $F_2^{\text{neq}}(T_a)$  and  $F_1^{\text{neq}}(T_b)$  defined in (32a) and (32d). The irreversible heats for the two isothermal processes are

$$\overline{Q}_{12}(T_a) = T_a \left[ S_1(T_a) - S_2(T_a) + \int_{\Omega} \rho_1(x; T_a) \ln \left( \frac{\rho_1(x; T_a)}{\rho_2(x; T_a)} \right) dx \right], \tag{38}$$

$$\overline{Q}_{21}(T_b) = T_b \left[ S_2(T_b) - S_1(T_b) + \int_{\Omega} \rho_2(x; T_b) \ln \left( \frac{\rho_2(x; T_b)}{\rho_1(x; T_b)} \right) dx \right]. \tag{39}$$

Meanwhile, those for the isochoric cooling and heating processes are

$$Q_2(T_a, T_b) = T_b \left[ S_2(T_a) - S_2(T_b) + \int_{\Omega} \rho_2(x; T_a) \ln \left( \frac{\rho_2(x; T_a)}{\rho_2(x; T_b)} \right) dx \right], \tag{40}$$

$$Q_1(T_b, T_a) = T_a \left[ S_1(T_b) - S_1(T_a) + \int_{\Omega} \rho_1(x; T_b) \ln \left( \frac{\rho_1(x; T_b)}{\rho_1(x; T_a)} \right) dx \right]. \tag{41}$$

Summarizing the whole heat cycle, we find that

$$\begin{aligned} W &= -(\overline{W}_{12} + \overline{W}_{21}) = -\overline{Q}_{12}(T_a) - Q_2(T_a, T_b) - \overline{Q}_{21}(T_b) - Q_1(T_b, T_a) \\ &\leq (T_a - T_b)[S_2(T_a) - S_1(T_b)]. \end{aligned} \tag{42}$$

This will lead to the same conclusions on the first-law and second-law efficiency for the Stirling cycle.

**Realization of a reversible cycle.** The Carnot cycle and Stirling cycle considered above are not truly reversible, once  $U_1 \neq U_2$  or  $T_a \neq T_b$ . To achieve the theoretical maximal efficiency, we need to construct a reversible heat cycle through a series of quasi-static processes, each of which involves only an infinitesimal change in either  $U$  or  $T$ . Taking the Stirling

cycle as an example. In the first isothermal working step, we insert  $N - 1$  intermediate states between  $\{T_a, U_1\}$  and  $\{T_a, U_2\}$ , that are  $\{T_a, U_1 + \Delta U\}, \{T_a, U_1 + 2\Delta U\}, \dots, \{T_a, U_1 + (N - 1)\Delta U\}$  with  $\Delta U = (U_2 - U_1)/N$ . In the limit of  $N \rightarrow \infty, \Delta U \rightarrow 0$ , which means each transition between two adjacent states can be treated as a quasi-static process. Therefore, the whole step between  $\{T_a, U_1\}$  and  $\{T_a, U_2\}$  becomes reversible with the help of those intermediate states. Applying similar procedure to other three steps, we will achieve a true thermodynamically reversible Stirling cycle by requiring an infinitesimal change in either  $U$  or  $T$  for each sub-step.

### 3.5. Work as a conditional expectation in energy representation

Consider once again two distributions  $\rho(x)$  and  $p^{eq}(x)$  with respective energy representations,  $\rho_E(y) = Z_1^{-1}(\beta_a)\Omega_1(y)e^{-\beta_a y}$  and  $p_E^{eq}(y) = Z_2^{-1}(\beta_b)\Omega_2(y)e^{-\beta_b y}$ . The key thermodynamic quantity that arises in (22), the irreversible work, cannot be expressed in terms of the six quantities:  $\Omega_1(y), Z_1(\beta), \Omega_2(y), Z_2(\beta)$ , and  $\beta_a, \beta_b$ . We note that

$$\int_{\Omega} \rho(x)[U_2(x) - U_1(x)]dx = \int_{\mathbb{R}} \left( \frac{\Omega_1(y)e^{-\beta_a y}}{Z_1(\beta_a)} \right) \{ \bar{U}_{2|U_1=y} - y \} dy, \tag{43}$$

in which

$$\bar{U}_{2|U_1=y} = \frac{\int_{y < U_1(x) \leq y+dh} U_2(x)dx}{\int_{y < U_1(x) \leq y+dh} dx}, \tag{44}$$

is a conditional expectation of  $U_2(x)$  given  $U_1(x) = y$ . The energy functions  $U_1(x)$  and  $U_2(x)$  are only two observables on the probability space and they certainly do not provide a full description of the probability space. Actually, knowing the canonical energy distributions  $\rho_E(y)$  and  $p_E^{eq}(y)$  is not equivalent to knowing their joint probability distribution; the missing information on their correlation is captured precisely in (44).

The lhs of (43) can also be expressed as

$$\begin{aligned} & \int_{\Omega} \rho(x)[U_2(x) - U_1(x)]dx \\ &= \frac{1}{\beta_a} \left[ \ln \frac{Z_1(\beta_a)}{Z_2(\beta_a)} + \int_{\Omega} \rho(x) \ln \left\{ \frac{e^{-\beta_a U_1(x)} Z_2(\beta_a)}{Z_1(\beta_a) e^{-\beta_a U_2(x)}} \right\} dx \right]. \end{aligned} \tag{45}$$

The term inside  $\{\cdot\cdot\}$  indeed can be understood as a Radon–Nikodym derivative between the two probability measures, which is well-defined on the entire  $\sigma$ -algebra  $\mathcal{F}$  as well as the restricted joint  $\sigma$ -algebra  $\mathcal{F}_{U_1, U_2}$ . However, it is singular on the further restricted  $\sigma$ -algebra  $\mathcal{F}_{U_1}$  or  $\mathcal{F}_{U_2}$ .

### 3.6. The role and consequence of determinism

Consider a sequence of measures  $\mu_{\epsilon}$  and two real-valued continuous random variables  $\mathbf{x}(\omega)$  and  $\mathbf{y}(\omega)$ , with corresponding probability density functions  $p_{\epsilon}(x)$  and  $q_{\epsilon}(x)$ . Their relative entropy is then

$$H[\mathbf{x} \parallel \mathbf{y}; \mu_{\epsilon}] = \int_{\Omega} p_{\epsilon}(x) \ln \left( \frac{p_{\epsilon}(x)}{q_{\epsilon}(x)} \right) dx. \tag{46}$$

If the sequence of measures  $\mu_{\epsilon}$  tends to a singleton with corresponding  $p_{\epsilon}(x) \rightarrow \delta(x - z)$  and  $q_{\epsilon}(x) \rightarrow \delta(x - y^*)$  as  $\epsilon \rightarrow 0$ , we call the limit *deterministic*.

It can be shown under rather weak conditions, or more properly through the theory of large deviations, that as  $\epsilon \rightarrow 0$  the  $p_{\epsilon}(x)$  and  $q_{\epsilon}(x)$  have asymptotic forms

$$\ln p_{\epsilon}(x) = -\frac{\varphi_p(x)}{\epsilon} + O(\ln \epsilon), \quad \ln q_{\epsilon}(x) = -\frac{\varphi_q(x)}{\epsilon} + O(\ln \epsilon), \tag{47}$$

in which  $\varphi_p(z) = \varphi_q(y^*) = 0$ . This asymptotic relation is known as the large deviations principle in the theory of probability [37]. Therefore,

$$\ln H[\mathbf{x} \parallel \mathbf{y}; \mu_{\epsilon}] \sim \frac{\varphi_q(z)}{\epsilon} + O(\ln \epsilon), \tag{48}$$

as  $\mathbf{x} \rightarrow z$ . Even though  $\mathbf{y} \rightarrow y^*$ , the relative entropy in (46) provides the  $\varphi_q$  as a function of  $z$  fully supported on  $\mathbb{R}^n$ . If the  $q_{\epsilon}$  is an invariant measure of a stochastic dynamical system, then the  $\varphi_q(z)$  is thought of as a “deterministic energy function”, which can be obtained as the asymptotic limit of determinism. The normalization of  $e^{-\varphi_q(x)/\epsilon}$ , however, is lost in the  $\ln \epsilon$ -order term. This corresponds to a certain gauge freedom.

A combination of the determinism with the canonical distribution immediately yields a key relationship that is well known in thermodynamics. Specifically, if the probability density function

$$\frac{\Omega^{(B)}(E)e^{-\beta E}}{Z(\beta)} = \frac{e^{-\beta E + \ln \Omega^{(B)}(E)}}{Z(\beta)} \rightarrow \delta(E - E^*), \tag{49}$$

in an asymptotic limit, then one has the *equation of state*

$$\left[ \frac{d}{dE} (\beta E - \ln \Omega^{(B)}(E)) \right]_{E=E^*} = 0. \tag{50}$$

A system in macroscopic thermodynamic equilibrium possesses one less degree of freedom [10]. Eq. (50) implies

$$\beta = \frac{d \ln \Omega^{(B)}(E^*)}{dE} = \frac{\frac{d}{dE} \Omega^{(B)}(E^*)}{\Omega^{(B)}(E^*)}, \tag{51}$$

in which

$$\begin{aligned} \Omega^{(B)}(E) &= \frac{1}{dE} \int_{E < U(x) \leq E + dE} dx = \oint_{U(x)=E} \frac{d\Sigma \cdot \hat{\mathbf{n}}}{\|\nabla U(x)\|} \\ &= \int_{U(x) \leq E} \nabla \cdot \left( \frac{\nabla U(x)}{\|\nabla U(x)\|^2} \right) dx, \end{aligned} \tag{52}$$

$$\begin{aligned} \frac{d\Omega^{(B)}(E)}{dE} &= \frac{1}{dE} \int_{E < U(x) \leq E + dE} \nabla \cdot \left( \frac{\nabla U(x)}{\|\nabla U(x)\|^2} \right) dx \\ &= \oint_{U(x) \leq E} \nabla \cdot \left( \frac{\nabla U(x)}{\|\nabla U(x)\|^2} \right) \frac{d\Sigma \cdot \hat{\mathbf{n}}}{\|\nabla U(x)\|}. \end{aligned} \tag{53}$$

Therefore,

$$\frac{d \ln \Omega^{(B)}(E)}{dE} = \frac{\oint_{U(x)=E} \nabla \cdot \left( \frac{\nabla U(x)}{\|\nabla U(x)\|^2} \right) \frac{d\Sigma \cdot \hat{\mathbf{n}}}{\|\nabla U(x)\|}}{\oint_{U(x)=E} \frac{d\Sigma \cdot \hat{\mathbf{n}}}{\|\nabla U(x)\|}}. \tag{54}$$

That is, the equilibrium  $\beta$  is the average of

$$\nabla \cdot \left( \frac{\nabla U}{\|\nabla U\|^2} \right) = \frac{\|\nabla U\| \nabla^2 U - 2 \nabla U \cdot \nabla \|\nabla U\|}{\|\nabla U\|^3}, \tag{55}$$

on the level-surface  $\{x : U(x) = E^*\}$ . For a given energy function  $U(x)$ , or an observable [14], Eq. (54), which generalizes the virial theorem in classical mechanics, provides the function  $\beta(E)$ .

### 4. The space of probability measures

#### 4.1. Affine structure, canonical distribution and its energy representation

We will now give a brief, non-rigorous introduction to the theory developed in [9]. Let  $\mathcal{M}$  be the set of all probability measures on  $(\Omega, \mathcal{F})$  that are absolutely continuous w.r.t. some probability measure  $\mathbb{P}$  (and therefore absolutely continuous w.r.t. each other) and let  $\mathcal{V}$  be an appropriate set of real-valued functions on  $\Omega$ . (Note that any choice of  $\mathbb{P}$  in  $\mathcal{M}$  would do; one only cares that all measures in  $\mathcal{M}$  are absolutely continuous w.r.t each other.) One now defines  $\oplus: \mathcal{M} \times \mathcal{V} \rightarrow \mathcal{M}$  such that

$$(\mu \oplus g)(A) = \frac{\int_A e^g d\mu}{\int_{\Omega} e^g d\mu}, \tag{56}$$

for any  $A \in \mathcal{F}$ . Assuming the denominator is finite (which requires some assumptions on  $\mathcal{V}$ ), the positivity of  $e^g$  implies that  $(\mu \oplus g)$  is also absolutely continuous w.r.t.  $\mathbb{P}$ . Since  $(\mu \oplus g)(\Omega) = 1$ , it is a probability measure. These two facts mean that  $(\mu \oplus g) \in \mathcal{M}$ , so the operation  $\oplus$  is well-defined. Note that  $\mu \oplus g = \mu \oplus (g + c)$  for any constant  $c$ , so this addition is not actually one-to-one. We can remedy this issue by restricting  $\mathcal{V}$  to functions that sum to zero, or we can replace each function with an equivalence class of functions that differ by a constant. One can then show that  $(\mathcal{M}, \mathcal{V}, \oplus)$  is an affine structure on  $\mathcal{M}$  [8,9]. If one chooses a particular measure  $\mathbb{P} \in \mathcal{M}$  as the origin, then any other measure  $\mu \in \mathcal{M}$  will have a Radon–Nikodym derivative  $\frac{d\mu}{d\mathbb{P}}(\omega)$ , and  $\mu = (\mathbb{P} \oplus g)$  where  $g = \ln\left(\frac{d\mu}{d\mathbb{P}}\right)$ .

Let  $J \subseteq \mathbb{R}$  be an interval and  $U \in \mathcal{V}$ . The function  $p: J \rightarrow \mathcal{M}$  such that  $p(\beta) = \mathbb{P} \oplus (-\beta U)$  is an affine straight line. More explicitly, we have the family of probability densities

$$\frac{e^{-\beta U(\omega)}}{Z(\beta)} \mathbb{P}(d\omega), \tag{57}$$

where  $Z(\beta)$  is the normalization factor.

In Kolmogorov's theory, the real-valued function  $U(\omega)$ , when thought of as a random variable, has its own probability density function w.r.t. the Lebesgue measure:

$$\mathbb{P}\{y < U(\omega) \leq y + dh\} = \frac{\Omega_U(y)}{Z(\beta)} e^{-\beta y} dh, \tag{58}$$

in which  $\ln \Omega_U(y)$  is the Gibbs entropy associated with function  $U(\omega)$ , defined in Eq. (12):

$$\Omega_U(y) = \frac{1}{dh} \int_{y < U(\omega) \leq y + dh} \mathbb{P}(d\omega). \tag{59}$$

The relation between the distributions in (57) and (58) establishes a map between the observables in the tangent space  $\mathcal{V}$  of  $\mathcal{M}$  and the standard probability density functions. (This is analogous to the dual relation between the Koopman operator on the space of observables and the Perron–Frobenius operator on the space of densities in dynamical systems theory.) We call (57) the *canonical representation* for the space of probability measures (SoPMs), and (58) its *energy representation*. Note that the energy representation of a given probability measure is not unique. The choice of  $U$  depends on both  $\mathbb{P}$  and  $\beta$ .

**A pair of observables.** We now discuss the notions of joint, marginal, and conditional probability in terms of the canonical representation in  $\mathcal{V}$ , with a fixed “origin”  $\mathbb{P}$ , which should be thought of as the  $\mathbb{P}$  in the probability space  $(\Omega, \mathcal{F}, \mathbb{P})$ , à la Kolmogorov. The SoPMs  $\mathcal{M}$  then is represented by observables  $U(\omega) \in \mathcal{V}$ , the tangent space of  $\mathcal{M}$ .

Consider two observables  $(U_1(\omega), U_2(\omega))$ , where  $U_2 \neq aU_1 + b$ . The corresponding “flat plane” can be parametrized as

$$\frac{e^{-\beta_a U_1(\omega) - \beta_b U_2(\omega)}}{Z_{1,2}(\beta_a, \beta_b)} \mathbb{P}(d\omega), \quad (\beta_a, \beta_b) \in \mathbb{R}^2. \tag{60}$$

We note that each observable induced a restricted  $\sigma$ -algebra on  $\mathbb{R}$ :  $\mathcal{F}_{U_1}$  and  $\mathcal{F}_{U_2}$  respectively, and the joint observable induces  $\mathcal{F}_{U_1, U_2} = \sigma(\mathcal{F}_{U_1} \cup \mathcal{F}_{U_2})$ . With respect to  $\mathcal{F}_{U_1, U_2}$ , the distribution in (60) can expressed on as:

$$\frac{\Omega_{1,2}(y_1, y_2)}{Z_{1,2}(\beta_a, \beta_b)} e^{-\beta_a y_1 - \beta_b y_2} dy_1 dy_2, \tag{61a}$$

in which

$$\Omega_{1,2}(y_1, y_2) = \frac{1}{dy_1 dy_2} \int_{y_1 < U_1(\omega) \leq y_1 + dy_1, y_2 < U_2(\omega) \leq y_2 + dy_2} d\mathbb{P}(\omega) \tag{61b}$$

$$= \frac{\partial^2}{\partial y_1 \partial y_2} \int_{U_1(\omega) \leq y_1, U_2(\omega) \leq y_2} d\mathbb{P}(\omega). \tag{61c}$$

We note that the marginal distribution

$$\int_{\mathbb{R}} \frac{\Omega_{1,2}(y_1, y_2)}{Z_{1,2}(\beta_a, \beta_b)} e^{-\beta_a y_1 - \beta_b y_2} dy_2 = \frac{1}{Z_1(\beta_a)} \left( \int_{\mathbb{R}} \Omega_{1,2}(y_1, y_2) dy_2 \right) e^{-\beta_a y_1}. \tag{62}$$

This implies that

$$\frac{1}{Z_1(\beta_a)} \int_{\mathbb{R}} \Omega_{1,2}(y_1, y_2) dy_2 = \frac{1}{Z_{1,2}(\beta_a, \beta_b)} \int_{\mathbb{R}} \Omega_{1,2}(y_1, y_2) e^{-\beta_b y_2} dy_2. \tag{63}$$

Since the rhs of (63) is not a function of  $\beta_b$ , we have the following equality:

$$\frac{\partial \ln Z_{1,2}(\beta_a, \beta_b)}{\partial \beta_b} = - \frac{\int_{\mathbb{R}} y_2 \Omega_{1,2}(y_1, y_2) e^{-\beta_b y_2} dy_2}{\int_{\mathbb{R}} \Omega_{1,2}(y_1, y_2) e^{-\beta_b y_2} dy_2}. \tag{64}$$

Eq. (63) can also be re-arranged into

$$\frac{\int_{\mathbb{R}} \Omega_{1,2}(y_1, y_2) e^{-\beta_b y_2} dy_2}{\int_{\mathbb{R}} \Omega_{1,2}(y_1, y_2) dy_2} = \frac{Z_{1,2}(\beta_a, \beta_b)}{Z_1(\beta_a)}, \tag{65a}$$

$$\frac{\int_{\mathbb{R}} \Omega_{1,2}(y_1, y_2) e^{-\beta_b y_2} \left( e^{\beta_b y_2} \right) dy_2}{\int_{\mathbb{R}} \Omega_{1,2}(y_1, y_2) e^{-\beta_b y_2} dy_2} = \frac{Z_1(\beta_a)}{Z_{1,2}(\beta_a, \beta_b)}. \tag{65b}$$

**Relative entropy between two random variables.** The relative entropy between two measures  $\mu_1 = (\mathbb{P} \oplus (-\beta_a U_1)) \in \mathcal{M}$  and  $\mu_2 = (\mathbb{P} \oplus (-\beta_b U_2(\omega))) \in \mathcal{M}$ , when transformed into the energy representation, is given by:

$$\begin{aligned} \int_{\Omega} \ln \left( \frac{d\mu_1}{d\mu_2}(\omega) \right) d\mu_1(\omega) &= \int_{\Omega} \frac{e^{-\beta_a U_1(\omega)}}{Z_1(\beta_a)} \ln \left( \frac{Z_2(\beta_b)}{Z_1(\beta_a)} e^{-\beta_a U_1(\omega) + \beta_b U_2(\omega)} \right) d\mathbb{P}(\omega) \\ &= \int_{\mathbb{R}} \frac{\Omega_{U_1}(h) e^{-\beta_a h}}{Z_1(\beta_a)} \ln \left( \frac{\Omega_{U_2}(h) e^{-\beta_a h}}{Z_1(\beta_a) \Omega_{U_1}(h)} \right) dh \end{aligned} \tag{66a}$$

$$+ \beta_b \int_{\Omega} U_2(\omega) \left( \frac{e^{-\beta_a U_1(\omega)}}{Z_1(\beta_a)} \right) d\mathbb{P}(\omega) + \ln Z_2(\beta_b). \tag{66b}$$

Note that the first term in (66b) again contains the  $\bar{U}_{2|U_1=h_1}$  that appeared in (25) and (44). It cannot be expressed in terms of the energy representations of  $\mu_1$  and  $\mu_2$ . Unless  $g_2 = ag_1 + b$ , the two measures  $\mu_1$  and  $\mu_2$ , with densities  $d\mu_1 = e^{g_1} d\mathbb{P}$  and  $d\mu_2 = e^{g_2} d\mathbb{P}$ , do not share the same restricted  $\sigma$ -algebra.

4.2. Entropy divergence in the SoPMs

Consider two probability measures  $\mu_1, \mu_2 \in \mathcal{M}$  in the SoPMs, with Radon–Nikodym derivatives w.r.t.  $\mathbb{P}$  given by  $f_1(\omega)$  and  $f_2(\omega)$  respectively. One can introduce the following divergence on  $\mathcal{M}$ :

$$d^2(\mu_1, \mu_2) = \int_{\Omega} (f_1(\omega) - f_2(\omega)) \left( \frac{\ln f_1(\omega) - \ln f_2(\omega)}{f_1(\omega) - f_2(\omega)} \right) (f_1(\omega) - f_2(\omega)) \mathbb{P}(d\omega). \tag{67}$$

This divergence can also be rewritten as the sum of two non-negative terms in the form of relative entropy, a symmetrized version of the latter:

$$d^2(\mu_1, \mu_2) = \int_{\Omega} \ln \left( \frac{d\mu_1}{d\mu_2}(\omega) \right) \mu_1(d\omega) + \int_{\Omega} \ln \left( \frac{d\mu_2}{d\mu_1}(\omega) \right) \mu_2(d\omega). \tag{68}$$

From this second form, it is clear that  $d$  is symmetric with respect to  $\mu_1$  and  $\mu_2$  and is zero if and only if  $\mu_1 = \mu_2$  on  $\mathcal{F}$ . This form also has the advantage of making it clear that  $d$  is invariant with respect to the choice of an origin  $\mathbb{P}$ . Note that, despite our notation, this quantity is not a metric because it does not satisfy the triangle inequality. It is only a local metric. That is, if  $\mu_1, \mu_2$  and  $\mu_3$  are sufficiently close together then  $d(\mu_1, \mu_2) + d(\mu_2, \mu_3) \geq d(\mu_1, \mu_3)$ .

**Divergence in energy representation.** If  $\mu_1$  and  $\mu_2$  are written in their respective energy representations, i.e.  $f_{E1}(y_1) = Z_1(\beta_a) \Omega_1(y_1) e^{-\beta_a y_1}$  and  $f_{E2}(y_2) = Z_2(\beta_b) \Omega_2(y_2) e^{-\beta_b y_2}$ . Then from Eq. (68), we have

$$\begin{aligned} d^2(\mu_1, \mu_2) &= \beta_b \int_{\mathbb{R}} \left( \frac{\Omega_1(y_1) e^{-\beta_a y_1}}{Z_1(\beta_a)} \right) \bar{U}_{2|U_1=y_1} dy_1 - \beta_a \bar{U}_1(\beta_a) - \beta_b \bar{U}_2(\beta_b) \\ &\quad + \beta_a \int_{\mathbb{R}} \left( \frac{\Omega_2(y_2) e^{-\beta_b y_2}}{Z_2(\beta_b)} \right) \bar{U}_{1|U_2=y_2} dy_2. \end{aligned} \tag{69}$$

There are three interesting special cases:

**Different  $\beta$ 's and same  $\Omega$ .** If  $\Omega_1(y) = \Omega_2(y) = \Omega(y)$ ,

$$d^2(\mu_1, \mu_2) = (\beta_b - \beta_a) (\bar{U}(\beta_a) - \bar{U}(\beta_b)). \tag{70}$$

**Different  $\Omega$ 's and same  $\beta$ .** With same  $\beta_a = \beta_b = \beta$  but different  $\Omega$ 's,

$$d^2(\mu_1, \mu_2) = \beta \int_{\mathbb{R}} \left( \frac{\Omega_1(y_1) e^{-\beta y_1}}{Z_1(\beta)} \right) [\bar{U}_{2|U_1=y_1} - y_1] dy_1 \tag{71a}$$

$$+ \beta \int_{\mathbb{R}} \left( \frac{\Omega_2(y_2) e^{-\beta y_2}}{Z_2(\beta)} \right) [\bar{U}_{1|U_2=y_2} - y_2] dy_2 \tag{71b}$$

$$= \beta (\bar{W}_{12}(\beta) + \bar{W}_{21}(\beta)). \tag{71c}$$

Here, following (22) and (23), we have identified the terms in (71a) and (71b) as  $\bar{W}_{12}(\beta)$  and  $\bar{W}_{21}(\beta)$ , respectively.

**Different  $\Omega$ 's and  $\beta$ 's.**

$$d^2(\mu_1, \mu_2) = \beta_b \bar{W}_{12}(\beta_a) + \beta_a \bar{W}_{21}(\beta_b) + (\beta_b - \beta_a) (\bar{U}_1(\beta_a) - \bar{U}_2(\beta_b)). \tag{72}$$

Eq. (72) implies an inequality that, being different from (35) and (36), is based on Massieu–Planck potential:

$$-\left(\frac{\overline{\mathcal{W}}_{12}}{T_b} + \frac{\overline{\mathcal{W}}_{21}}{T_a}\right) \leq (\overline{U}_1 - \overline{U}_2) \left(\frac{1}{T_b} - \frac{1}{T_a}\right). \tag{73}$$

### 4.3. Heat divergence

One can also introduce another related divergence on  $\mathcal{M}$ . For fixed  $\beta_a, \beta_b > 0$ , define:

$$\begin{aligned} d_\beta^2(\mu_1, \mu_2) &= \frac{1}{\beta_a} \int_\Omega f_1(\omega) \ln \left(\frac{f_1(\omega)}{f_2^{(\beta_a)}(\omega)}\right) \mathbb{P}(d\omega) + \frac{1}{\beta_b} \int_\Omega f_2(\omega) \ln \left(\frac{f_2(\omega)}{f_1^{(\beta_b)}(\omega)}\right) \mathbb{P}(d\omega) \\ &= \int_\Omega \left(\frac{e^{-\beta_a U_1(\omega)}}{Z_1(\beta_a)} - \frac{e^{-\beta_b U_2(\omega)}}{Z_2(\beta_b)}\right) (U_2(\omega) - U_1(\omega)) \mathbb{P}(d\omega) \\ &\quad + \frac{1}{\beta_a} \ln \left(\frac{Z_2(\beta_a)}{Z_1(\beta_a)}\right) - \frac{1}{\beta_b} \ln \left(\frac{Z_2(\beta_b)}{Z_1(\beta_b)}\right), \end{aligned} \tag{74}$$

in which

$$f_1(\omega) = \frac{e^{-\beta_a U_1(\omega)}}{Z_1(\beta_a)} \text{ and } f_2(\omega) = \frac{e^{-\beta_b U_2(\omega)}}{Z_2(\beta_b)} \tag{75}$$

are the densities of  $\mu_1$  and  $\mu_2$  with respect to  $\mathbb{P}$  and

$$f_2^{(\beta_a)}(\omega) = \frac{e^{-\beta_a U_2(\omega)}}{Z_2(\beta_a)}, \quad f_1^{(\beta_b)}(\omega) = \frac{e^{-\beta_b U_1(\omega)}}{Z_1(\beta_b)}. \tag{76}$$

The same caveats as before apply: This is not a metric on  $\mathcal{M}$  because it does not satisfy the triangle inequality, but it is a local metric in the sense that the triangle inequality is satisfied when all measures are sufficiently close together. We shall call  $d_\beta(\cdot, \cdot)$  in (74) the *heat divergence*. In terms of

$$\mathcal{W}_{12}(\omega) = \frac{1}{\beta_a} \ln \left(\frac{e^{-\beta_a U_1(\omega)}}{e^{-\beta_a U_2(\omega)}}\right), \quad \mathcal{W}_{21}(\omega) = \frac{1}{\beta_b} \ln \left(\frac{e^{-\beta_b U_2(\omega)}}{e^{-\beta_b U_1(\omega)}}\right), \tag{77}$$

we have

$$\begin{aligned} d_\beta^2(\mu_1, \mu_2) &= \mathbb{E}^{\mu_1}[\mathcal{W}_{12}(\omega)] + \beta_a^{-1} \ln \mathbb{E}^{\mu_1}[e^{-\beta_a \mathcal{W}_{12}(\omega)}] + \mathbb{E}^{\mu_2}[\mathcal{W}_{21}(\omega)] \\ &\quad + \beta_b^{-1} \ln \mathbb{E}^{\mu_2}[e^{-\beta_b \mathcal{W}_{21}(\omega)}]. \end{aligned} \tag{78}$$

Using the Jarzynski–Crooks relation from (28), Eq. (78) implies

$$\overline{\mathcal{W}}_{12}(\beta_a) + \overline{\mathcal{W}}_{21}(\beta_b) + F_1(\beta_a) - F_2(\beta_a) + F_2(\beta_b) - F_1(\beta_b) \geq 0. \tag{79}$$

This result generalizes Carnot’s inequality.

### 4.4. Infinitesimal entropy metric associated with $\Delta\beta$

Consider an infinitesimal change in  $\beta \rightarrow \beta + \Delta\beta$  and corresponding  $d\mu = e^{-\beta U} d\mathbb{P} \rightarrow d(\mu + \Delta\mu) = e^{-(\beta + \Delta\beta)U} d\mathbb{P}$ . Then we have

$$\begin{aligned} d^2(\mu, \mu + \Delta\mu) &= (\Delta\beta)^2 \int_0^\infty \frac{\Omega(y)e^{-\beta y}}{Z(\beta)} \left[\left(\frac{d \ln Z}{d\beta}\right) + y\right]^2 dy \\ &= (\Delta\beta)^2 \int_0^\infty \frac{\Omega(y)e^{-\beta y}}{Z(\beta)} (y - \mathbb{E}[U])^2 dy \\ &= (\Delta\beta)^2 \text{Var}[U]. \end{aligned} \tag{80}$$

This is a very important relation that connects the *entropy divergence* with *temperature* and *energy fluctuations*. Furthermore, we have

$$d^2(\mu, \mu + \Delta\mu) = (\Delta\beta)^2 \left(-\frac{d^2 \ln Z(\beta)}{d\beta^2}\right) = (\Delta\beta)^2 \left(\frac{d}{d\beta} \mathbb{E}[U]\right). \tag{81}$$

The term inside  $(\dots)$  on the rhs is called the *heat capacity* in thermodynamics. Internal energy  $\mathbb{E}[U]$  is a “slope” and the  $\text{Var}[X_\beta]$  is a curvature of the “potential function”  $-\ln Z(\beta)$ .

4.5. A mathematical remark

**Log-mean-exponential inequality and equality.** We see that both entropy divergence in (68) and heat divergence in (78) are based on a very general inequality involving the log-mean-exponential of a random variable  $\xi(\omega)$  [38]: Jensen’s inequality.

$$\mathbb{E}[\xi(\omega)] + \beta^{-1} \ln \mathbb{E}[e^{-\beta\xi(\omega)}] \geq 0. \tag{82}$$

In (68), the two  $\xi$ s are the information  $\ln \frac{d\mu_1}{d\mu_2}(\omega)$  and  $\ln \frac{d\mu_2}{d\mu_1}(\omega)$ ; and in (78), the two  $\xi$ s are the work  $\mathcal{W}_{12}(\omega) = \beta_a^{-1} \ln \frac{e^{-\beta_a U_1(\omega)}}{e^{-\beta_a U_2(\omega)}}$  and  $\mathcal{W}_{21}(\omega) = \beta_b^{-1} \ln \frac{e^{-\beta_b U_2(\omega)}}{e^{-\beta_b U_1(\omega)}}$ . They are all different forms of Radon–Nikodym derivatives. In the entropy divergence, the second, log-mean-exponential term in (82) is zero according to the Hatano–Sasa equality. In the heat divergence case, the same term gives a Jarzynski–Crooks’ free energy difference.

Eq. (82) should be recognized as “mean internal energy minus free energy”. Thus it should be some kind of entropy:

$$\mathbb{E}^{\mathbb{P}'}[\xi(\omega)] + \beta^{-1} \ln \mathbb{E}^{\mathbb{P}}[e^{-\beta\xi(\omega)}] = \mathbb{E}^{\mathbb{P}}\left[\ln\left(\frac{d\mathbb{P}}{d\mathbb{P}'}(\omega)\right)\right], \tag{83}$$

in which  $\mathbb{P}' = \mathbb{P} \oplus (-\beta\xi)$  is the affine sum of  $\mathbb{P}$  and  $(-\beta\xi)$ . Eq. (83) could be argued as the *fundamental equation for isothermal processes* under a single temperature  $T = \beta^{-1}$ . The implication of this interesting “Jensen’s equality” to the affine geometry of the SoPMs is currently being explored.

5. Discussion

It has been well established, through the work of Gibbs, Carathéodory, and many others, that geometry has a role in the theory of equilibrium thermodynamics [39–41]. Classical thermodynamics is not based on the theory of chance, but there is no doubt that the notion of entropy has its root in the theory of probability. In the present work, we propose that the space of probability measures as a natural setting in which thermodynamic concepts can be established logically. In particular, an affine structure is naturally related to the canonical probability distribution studied by Boltzmann and Gibbs in their statistical theories, and almost all thermodynamic potentials are different forms of Radon–Nikodym derivatives associated with *changes of measures*. Even the fundamental equation of nonequilibrium thermodynamics, together with the distinctly nonequilibrium notion of entropy production, naturally emerges.

Statistical mechanics, as a scientific theory, differs from Kolmogorov’s axiomatic theory of probability in one essential point: The latter demands a complete probability space and a normalized probability measure, while in the former every probability distribution is a *conditioned probability* under many known and unknown conditions. More importantly, the probability of the conditions, themselves as random events, are usually not knowable. In the theory of the space of measures, we see that one mechanical system with a given energy function  $U(\omega)$  corresponds to a straight line, and the fixing of the origin in  $\mathcal{M}$  in terms of  $\mathbb{P}$  or the normalization in terms of  $Z(\beta)$  [which translates to the arbitrary constant in  $U(\omega)$ ] amount to the idea of gauge fixing. Thermodynamic work then arises in the rotation from  $U_1(\omega)$  to  $U_2(\omega)$ . In the theory of probability, associated with any “change” is a *change of measure*: Radon–Nikodym derivatives simply provide the calculus to quantify the *fluxion*! In Newtonian mechanics, change in space is absolute; but in probability, it is a complex matter, and it is all relative.

The probability theory of large deviations is now a recognized mathematical foundation for statistical thermodynamics [37,42,43]. Such a theory is concerned with the deterministic *thermodynamic limit*. In Section 3.6, we see that the combination of our theory and a deterministic limit gives rise to the concept of *macroscopic equations of state* in classic thermodynamics [10].

Equilibrium mean internal energy  $\bar{U}(\beta)$  depends on both the intrinsic properties of a system and its external environment. This is most clearly shown through the canonical distribution that is determined by  $U(\omega)$  and  $\beta$ . The decomposition in Eq. (15), a simple example of the much more general (83), connects the internal energy with “work” and “heat”, or the “usable energy” and “useless energy”, or entropy production and entropy change. These are all just different interpretations under different perspectives.

Acknowledgments

We thank Yu-Chen Cheng and Ying-Jen Yang for many helpful discussions, and Professors Jin Feng (University of Kansas) and Hao Ge (Peking University) for advices. L.H. acknowledges the financial supports from the National Natural Science Foundation of China (Grants 21877070) and Tsinghua University Initiative Scientific Research Program (Grants 20151080424). H.Q. acknowledges the Olga Jung Wan Endowed Professorship for support.

## References

- [1] A.N. Kolmogorov, S.V. Fomin, *Introductory Real Analysis*, Silverman, R.A. Transl., Dover, New York, 1968.
- [2] H. Qian, Mesoscopic nonequilibrium thermodynamics of single macromolecules and dynamic entropy-energy compensation, *Phys. Rev. E* 65 (2001) 016102.
- [3] F.X.F. Ye, H. Qian, *Stochastic dynamics II: Finite random dynamical systems, linear representation, and entropy production*, *Discrete Contin. Dyn. Syst. B* 24 (2019) 4341–4366.
- [4] W. Feller, The general diffusion operator and positive preserving semi-group in one dimension, *Ann. of Math.* 60 (1954) 417–436.
- [5] E. Nelson, An existence theorem for second order parabolic equations, *Trans. Amer. Math. Soc.* 88 (1958) 414–429.
- [6] D.Q. Jiang, M. Qian, M.P. Qian, *Mathematical Theory of Nonequilibrium Steady States*, Springer, New York, 2004.
- [7] G.A. Pavliotis, *Stochastic Processes and Applications*, Springer, New York, 2014.
- [8] J. Gallier, *Geometric Methods and Applications*, second ed., Springer, New York, 2011.
- [9] L.F. Thompson, *Affine structures, geometry and thermodynamics*, 2020, Manuscript in preparation.
- [10] W. Pauli, *Pauli Lectures on Physics: Vol. 3, Thermodynamics and the Kinetic Theory of Gases; Vol. 4, Statistical Mechanics*, The MIT Press, Cambridge, MA, 1973.
- [11] H. Qian, Y.-C. Cheng, L.F. Thompson, Ternary representation of stochastic change and the origin of entropy and its fluctuations, 2019, arXiv:1902.09536.
- [12] T.M. Cover, J.A. Thomas, *Elements of Information Theory*, John Wiley & Sons, New York, 1991.
- [13] Y.-J. Yang, H. Qian, Unified formalism for entropy production and fluctuation relations, *Phys. Rev. E* 101 (2020) 022129.
- [14] Y.-C. Cheng, H. Qian, Y. Zhu, Asymptotic behavior of a sequence of conditional probability distributions and the canonical ensemble, 2020, arXiv:1912.11137.
- [15] H. Qian, Information and entropic force: physical description of biological cells, chemical reaction kinetics, and information theory, *Sci. Sin. Vitae* 47 (2017) 257–261 (in Chinese).
- [16] A.N. Kolmogorov, Three approaches to the quantitative definition of information, *Int. J. Comput. Math.* 2 (1968) 157–168.
- [17] M. Tribus, *Thermostatistics and Thermodynamics: An Introduction to Energy, Information and States of Matter, with Engineering Applications*, D. van Nostrand, New York, 1961.
- [18] M.C. Mackey, The dynamic origin of increasing entropy, *Rev. Modern Phys.* 61 (1989) 981–1015.
- [19] H. Qian, Relative entropy: Free energy associated with equilibrium fluctuations and nonequilibrium deviations, *Phys. Rev. E* 63 (2001) 042103.
- [20] I. Prigogine, *Introduction to Thermodynamics of Irreversible Processes*, Charles C. Thomas Pub., Springfield, IL, 1955.
- [21] H. Qian, Thermodynamics of the general diffusion process: Equilibrium supercurrent and nonequilibrium driven circulation with dissipation, *Eur. Phys. J. Spec. Top.* 224 (2015) 781–799.
- [22] H. Qian, S. Kjelstrup, A.B. Kolomeisky, D. Bedeaux, Entropy production in mesoscopic stochastic thermodynamics - Nonequilibrium kinetic cycles driven by chemical potentials, temperatures, and mechanical forces, *J. Phys. Condens. Matter.* 28 (2016) 153004.
- [23] R.W. Zwanzig, High-temperature equation of state by a perturbation method. I. Nonpolar gases, *J. Chem. Phys.* 22 (1954) 1420–1426.
- [24] C. Jarzynski, Nonequilibrium equality for free energy differences, *Phys. Rev. Lett.* 78 (1997) 2690–2693.
- [25] G. Crooks, Entropy production fluctuation theorem and the nonequilibrium work relation for free energy differences, *Phys. Rev. E* 60 (1999) 2721–2726.
- [26] T. Hatano, S.I. Sasa, Steady-state thermodynamics of Langevin systems, *Phys. Rev. Lett.* 86 (2001) 3463–3466.
- [27] P. Ao, H. Qian, Y. Tu, J. Wang, A theory of mesoscopic phenomena: Time scales, emergent unpredictability, symmetry breaking and dynamics across different levels, 2013, arXiv:1310.5585.
- [28] R. Balian, *From Microphysics to Macrophysics: Methods and Applications of Statistical Physics*, Vol. I, Springer, New York, 1991.
- [29] D. Frenkel, P.B. Warren, Gibbs, Boltzmann, and negative temperatures, *Amer. J. Phys.* 83 (2015) 163–170.
- [30] R.M. Noyes, Entropy of mixing of interconvertible species: Some reflections on the Gibbs paradox, *J. Chem. Phys.* 34 (1961) 1983–1985.
- [31] J.G. Kirkwood, Statistical mechanics of fluid mixtures, *J. Chem. Phys.* 3 (1935) 300–313.
- [32] A. Ben-Naim, Mixing and assimilation in systems of interaction particles, *Amer. J. Phys.* 55 (1987) 1105–1109.
- [33] J.M.R. Parrondo, J.M. Horowitz, T. Sagawa, Thermodynamics of information, *Nat. Phys.* 11 (2015) 131–139.
- [34] T.M. Hoang, R. Pan, J. Ahn, J. Bang, H.T. Quan, T. Li, Experimental test of the differential fluctuation theorem and a generalized Jarzynski equality for arbitrary initial states, *Phys. Rev. Lett.* 120 (2018) 080602.
- [35] J. Honerkamp, *Statistical Physics: An Advanced Approach with Applications*, Springer, New York, 2002.
- [36] L. Chen, C. Wu, F. Sun, Finite time thermodynamic optimization or entropy generation minimization of energy systems, *J. Non-Equilib. Thermodyn.* 24 (1999) 327–359.
- [37] H. Touchette, The large deviation approach to statistical mechanics, *Phys. Rep.* 478 (2009) 1–69.
- [38] H. Qian, Nonequilibrium potential function of chemically driven single macromolecules via Jarzynski-type log-mean-exponential heat, *J. Phys. Chem. B* 109 (2005) 23624–23628.
- [39] G.A. Maugin, *Continuum Mechanics Through the Eighteenth and Nineteenth Centuries: Historical Perspectives from Bernoulli to Hellinger*, Springer, New York, 2014, pp. 137–147.
- [40] L. Pogliani, M.N. Berberan-Santos, Constantin Carathéodory and the axiomatic thermodynamics, *J. Math. Chem.* 28 (2000) 313–324.
- [41] P. Salamon, B. Andresen, J. Nulton, A.K. Konopka, The mathematical structure of thermodynamics, in: A.K. Konopka (Ed.), *Handbook of Systems Biology*, CRC Press, Boca Raton, 2006.
- [42] R.S. Ellis, *Entropy, Large Deviations, and Statistical Mechanics*, Springer, New York, 2006.
- [43] H. Ge, H. Qian, Mesoscopic kinetic basis of macroscopic chemical thermodynamics: A mathematical theory, *Phys. Rev. E* 94 (2016) 052150.