

# Algorithmic Thermodynamics

John C. Baez

Department of Mathematics, University of California  
Riverside, California 92521, USA

Mike Stay

Computer Science Department, University of Auckland  
*and*

Google, 1600 Amphitheatre Pkwy  
Mountain View, California 94043, USA

email: baez@math.ucr.edu, stay@google.com

September 5, 2015

## Abstract

Algorithmic entropy can be seen as a special case of entropy as studied in statistical mechanics. This viewpoint allows us to apply many techniques developed for use in thermodynamics to the subject of algorithmic information theory. In particular, suppose we fix a universal prefix-free Turing machine and let  $X$  be the set of programs that halt for this machine. Then we can regard  $X$  as a set of ‘microstates’, and treat any function on  $X$  as an ‘observable’. For any collection of observables, we can study the Gibbs ensemble that maximizes entropy subject to constraints on expected values of these observables. We illustrate this by taking the log runtime, length, and output of a program as observables analogous to the energy  $E$ , volume  $V$  and number of molecules  $N$  in a container of gas. The conjugate variables of these observables allow us to define quantities which we call the ‘algorithmic temperature’  $T$ , ‘algorithmic pressure’  $P$  and ‘algorithmic potential’  $\mu$ , since they are analogous to the temperature, pressure and chemical potential. We derive an analogue of the fundamental thermodynamic relation  $dE = TdS - PdV + \mu dN$ , and use it to study thermodynamic cycles analogous to those for heat engines. We also investigate the values of  $T$ ,  $P$  and  $\mu$  for which the partition function converges. At some points on the boundary of this domain of convergence, the partition function becomes uncomputable. Indeed, at these points the partition function itself has nontrivial algorithmic entropy.

# 1 Introduction

Many authors [1, 6, 9, 12, 16, 24, 26, 28] have discussed the analogy between algorithmic entropy and entropy as defined in statistical mechanics: that is, the entropy of a probability measure  $p$  on a set  $X$ . It is perhaps insufficiently appreciated that algorithmic entropy can be seen as a *special case* of the entropy as defined in statistical mechanics. We describe how to do this in Section 3.

This allows all the basic techniques of thermodynamics to be imported to algorithmic information theory. The key idea is to take  $X$  to be some version of ‘the set of all programs that eventually halt and output a natural number’, and let  $p$  be a Gibbs ensemble on  $X$ . A Gibbs ensemble is a probability measure that maximizes entropy subject to constraints on the mean values of some observables — that is, real-valued functions on  $X$ .

In most traditional work on algorithmic entropy, the relevant observable is the length of the program. However, much of the interesting structure of thermodynamics only becomes visible when we consider several observables. When  $X$  is the set of programs that halt and output a natural number, some other important observables include the output of the program and logarithm of its runtime. So, in Section 4 we illustrate how ideas from thermodynamics can be applied to algorithmic information theory using these three observables.

To do this, we consider a Gibbs ensemble of programs which maximizes entropy subject to constraints on:

- $E$ , the expected value of the logarithm of the program’s runtime (which we treat as analogous to the energy of a container of gas),
- $V$ , the expected value of the length of the program (analogous to the volume of the container), and
- $N$ , the expected value of the program’s output (analogous to the number of molecules in the gas).

This measure is of the form

$$p = \frac{1}{Z} e^{-\beta E(x) - \gamma V(x) - \delta N(x)}$$

for certain numbers  $\beta, \gamma, \delta$ , where the normalizing factor

$$Z = \sum_{x \in X} e^{-\beta E(x) - \gamma V(x) - \delta N(x)}$$

is called the ‘partition function’ of the ensemble. The partition function reduces to Chaitin’s number  $\Omega$  when  $\beta = 0$ ,  $\gamma = \ln 2$  and  $\delta = 0$ . This number is uncomputable [6]. However, we show that the partition function  $Z$  is computable when  $\beta > 0$ ,  $\gamma \geq \ln 2$ , and  $\delta \geq 0$ .

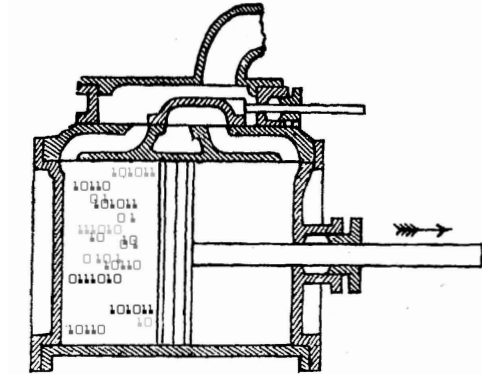
We derive an algorithmic analogue of the basic thermodynamic relation

$$dE = TdS - PdV + \mu dN.$$

Here:

- $S$  is the entropy of the Gibbs ensemble,
- $T = 1/\beta$  is the ‘algorithmic temperature’ (analogous to the temperature of a container of gas). Roughly speaking, this counts how many times you must double the runtime in order to double the number of programs in the ensemble while holding their mean length and output fixed.
- $P = \gamma/\beta$  is the ‘algorithmic pressure’ (analogous to pressure). This measures the tradeoff between runtime and length. Roughly speaking, it counts how much you need to decrease the mean length to increase the mean log runtime by a specified amount, while holding the number of programs in the ensemble and their mean output fixed.
- $\mu = -\delta/\beta$  is the ‘algorithmic potential’ (analogous to chemical potential). Roughly speaking, this counts how much the mean log runtime increases when you increase the mean output while holding the number of programs in the ensemble and their mean length fixed.

Starting from this relation, we derive analogues of Maxwell’s relations and consider thermodynamic cycles such as the Carnot cycle or Stoddard cycle. For this we must introduce concepts of ‘algorithmic heat’ and ‘algorithmic work’.



Charles Babbage described a computer powered by a steam engine; we describe a heat engine powered by programs! We admit that the significance of this line of thinking remains a bit mysterious. However, we hope it points the way toward a further synthesis of algorithmic information theory and thermodynamics. We call this hoped-for synthesis ‘algorithmic thermodynamics’.

## 2 Related Work

Li and Vitányi use the term ‘algorithmic thermodynamics’ for describing physical states using a universal prefix-free Turing machine  $U$ . They look at the

smallest program  $p$  that outputs a description  $x$  of a particular microstate to some accuracy, and define the physical entropy to be

$$S_A(x) = (k \ln 2)(K(x) + H_x),$$

where  $K(x) = |p|$  and  $H_x$  embodies the uncertainty in the actual state given  $x$ . They summarize their own work and subsequent work by others in chapter eight of their book [17]. Whereas they consider  $x = U(p)$  to be a microstate, we consider  $p$  to be the microstate and  $x$  the value of the observable  $U$ . Then their observables  $O(x)$  become observables of the form  $O(U(p))$  in our model.

Tadaki [27] generalized Chaitin’s number  $\Omega$  to a function  $\Omega^D$  and showed that the value of this function is compressible by a factor of exactly  $D$  when  $D$  is computable. Calude and Stay [5] pointed out that this generalization was formally equivalent to the partition function of a statistical mechanical system where temperature played the role of the compressibility factor, and studied various observables of such a system. Tadaki [28] then explicitly constructed a system with that partition function: given a total length  $E$  and number of programs  $N$ , the entropy of the system is the log of the number of  $E$ -bit strings in  $\text{dom}(U)^N$ . The temperature is

$$\frac{1}{T} = \left. \frac{\Delta E}{\Delta S} \right|_N.$$

In a follow-up paper [29], Tadaki showed that various other quantities like the free energy shared the same compressibility properties as  $\Omega^D$ . In this paper, we consider multiple variables, which is necessary for thermodynamic cycles, chemical reactions, and so forth.

Manin and Marcolli [20] derived similar results in a broader context and studied phase transitions in those systems. Manin [18, 19] also outlined an ambitious program to treat the infinite runtimes one finds in undecidable problems as singularities to be removed through the process of renormalization. In a manner reminiscent of hunting for the proper definition of the “one-element field”  $F_{un}$ , he collected ideas from many different places and considered how they all touch on this central theme. While he mentioned a runtime cutoff as being analogous to an energy cutoff, the renormalizations he presented are uncomputable. In this paper, we take the log of the runtime as being analogous to the energy; the randomness described by Chaitin and Tadaki then arises as the infinite-temperature limit.

### 3 Algorithmic Entropy

To see algorithmic entropy as a special case of the entropy of a probability measure, it is useful to follow Solomonoff [24] and take a Bayesian viewpoint. In Bayesian probability theory, we always start with a probability measure called a ‘prior’, which describes our assumptions about the situation at hand before we make any further observations. As we learn more, we may update this

prior. This approach suggests that we should define the entropy of a probability measure *relative to another probability measure* — the prior.

A probability measure  $p$  on a finite set  $X$  is simply a function  $p: X \rightarrow [0, 1]$  whose values sum to 1, and its entropy is defined as follows:

$$S(p) = - \sum_{x \in X} p(x) \ln p(x).$$

But we can also define the entropy of  $p$  relative to another probability measure  $q$ :

$$S(p, q) = - \sum_{x \in X} p(x) \ln \frac{p(x)}{q(x)}.$$

This **relative entropy** has been extensively studied and goes by various other names, including ‘Kullback–Leibler divergence’ [13] and ‘information gain’ [23].

The term ‘information gain’ is nicely descriptive. Suppose we initially assume the outcome of an experiment is distributed according to the probability measure  $q$ . Suppose we then repeatedly do the experiment and discover its outcome is distributed according to the measure  $p$ . Then the information gained is  $S(p, q)$ .

Why? We can see this in terms of coding. Suppose  $X$  is a finite set of signals which are randomly emitted by some source. Suppose we wish to encode these signals as efficiently as possible in the form of bit strings. Suppose the source emits the signal  $x$  with probability  $p(x)$ , but we erroneously believe it is emitted with probability  $q(x)$ . Then  $S(p, q)/\ln 2$  is the expected extra message-length per signal that is required if we use a code that is optimal for the measure  $q$  instead of a code that is optimal for the true measure,  $p$ .

The ordinary entropy  $S(p)$  is, up to a constant, just the relative entropy in the special case where the prior assigns an equal probability to each outcome. In other words:

$$S(p) = S(p, q_0) + S(q_0)$$

when  $q_0$  is the so-called ‘uninformative prior’, with  $q_0(x) = 1/|X|$  for all  $x \in X$ .

We can also define relative entropy when the set  $X$  is countably infinite. As before, a probability measure on  $X$  is a function  $p: X \rightarrow [0, 1]$  whose values sum to 1. And as before, if  $p$  and  $q$  are two probability measures on  $X$ , the entropy of  $p$  relative to  $q$  is defined by

$$S(p, q) = - \sum_{x \in X} p(x) \ln \frac{p(x)}{q(x)}. \tag{1}$$

But now the role of the prior becomes more clear, because there is no probability measure that assigns the same value to each outcome!

In what follows we will take  $X$  to be — roughly speaking — the set of all programs that eventually halt and output a natural number. As we shall see, while this set is countably infinite, there are still some natural probability measures on it, which we may take as priors.

To make this precise, we recall the concept of a universal prefix-free Turing machine. In what follows we use **string** to mean a bit string, that is, a finite, possibly empty, list of 0's and 1's. If  $x$  and  $y$  are strings, let  $x||y$  be the concatenation of  $x$  and  $y$ . A **prefix** of a string  $z$  is a substring beginning with the first letter, that is, a string  $x$  such that  $z = x||y$  for some  $y$ . A **prefix-free** set of strings is one in which no element is a prefix of any other. The **domain**  $\text{dom}(M)$  of a Turing machine  $M$  is the set of strings that cause  $M$  to eventually halt. We call the strings in  $\text{dom}(M)$  **programs**. We assume that when the  $M$  halts on the program  $x$ , it outputs a natural number  $M(x)$ . Thus we may think of the machine  $M$  as giving a function  $M: \text{dom}(M) \rightarrow \mathbb{N}$ .

A **prefix-free Turing machine** is one whose halting programs form a prefix-free set. A prefix-free machine  $U$  is **universal** if for any prefix-free Turing machine  $M$  there exists a constant  $c$  such that for each string  $x$ , there exists a string  $y$  with

$$U(y) = M(x) \text{ and } |y| < |x| + c.$$

Let  $U$  be a universal prefix-free Turing machine. Then we can define some probability measures on  $X = \text{dom}(U)$  as follows. Let

$$|\cdot|: X \rightarrow \mathbb{N}$$

be the function assigning to each bit string its length. Then there is for any constant  $\gamma > \ln 2$  a probability measure  $p$  given by

$$p(x) = \frac{1}{Z} e^{-\gamma|x|}.$$

Here the normalization constant  $Z$  is chosen to make the numbers  $p(x)$  sum to 1:

$$Z = \sum_{x \in X} e^{-\gamma|x|}.$$

It is worth noting that for computable real numbers  $\gamma \geq \ln 2$ , the normalization constant  $Z$  is uncomputable [27]. Indeed, when  $\gamma = \ln 2$ ,  $Z$  is Chaitin's famous number  $\Omega$ . We return to this issue in Section 4.5.

Let us assume that each program prints out some natural number as its output. Thus we have a function

$$N: X \rightarrow \mathbb{N}$$

where  $N(x)$  equals  $i$  when program  $x$  prints out the number  $i$ . We may use this function to 'push forward'  $p$  to a probability measure  $q$  on the set  $\mathbb{N}$ . Explicitly:

$$q(i) = \sum_{x \in X: N(x)=i} e^{-\gamma|x|}.$$

In other words, if  $i$  is some natural number,  $q(i)$  is the probability that a program randomly chosen according to the measure  $p$  will print out this number.

Given any natural number  $n$ , there is a probability measure  $\delta_n$  on  $\mathbb{N}$  that assigns probability 1 to this number:

$$\delta_n(m) = \begin{cases} 1 & \text{if } m = n \\ 0 & \text{otherwise.} \end{cases}$$

We can compute the entropy of  $\delta_n$  relative to  $q$ :

$$\begin{aligned} S(\delta_n, q) &= - \sum_{i \in \mathbb{N}} \delta_n(i) \ln \frac{\delta_n(i)}{q(i)} \\ &= - \ln \left( \sum_{x \in X: N(x)=n} e^{-\gamma|x|} \right) + \ln Z. \end{aligned} \tag{2}$$

Since the quantity  $\ln Z$  is independent of the number  $n$ , and uncomputable, it makes sense to focus attention on the other part of the relative entropy:

$$- \ln \left( \sum_{x \in X: N(x)=n} e^{-\gamma|x|} \right).$$

If we take  $\gamma = \ln 2$ , this is precisely the **algorithmic entropy** [7, 16] of the number  $n$ . So, up to the additive constant  $\ln Z$ , we have seen that *algorithmic entropy is a special case of relative entropy*.

One way to think about entropy is as a measure of surprise: if you can predict what comes next — that is, if you have a program that can compute it for you — then you are not surprised. For example, the first 2000 bits of the binary fraction for  $1/3$  can be produced with this short Python program:

```
print "01" * 1000
```

But if the number is complicated, if every bit is surprising and unpredictable, then the shortest program to print the number does not do any computation at all! It just looks something like

```
print "101000011001010010100101000101111101101101001010"
```

Levin's coding theorem [15] says that the difference between the algorithmic entropy of a number and its **Kolmogorov complexity** — the length of the shortest program that outputs it — is bounded by a constant that only depends on the programming language.

So, up to some error bounded by a constant, *algorithmic information is information gain*. The algorithmic entropy is the information gained upon learning a number, if our prior assumption was that this number is the output of a randomly chosen program — randomly chosen according to the measure  $p$  where  $\gamma = \ln 2$ .

So, algorithmic entropy is not just *analogous* to entropy as defined in statistical mechanics: it is a *special case*, as long as we take seriously the Bayesian

philosophy that entropy should be understood as relative entropy. This realization opens up the possibility of taking many familiar concepts from thermodynamics, expressed in the language of statistical mechanics, and finding their counterparts in the realm of algorithmic information theory.

But to proceed, we must also understand more precisely the role of the measure  $p$ . In the next section, we shall see that this type of measure is already familiar in statistical mechanics: it is a Gibbs ensemble.

## 4 Algorithmic Thermodynamics

Suppose we have a countable set  $X$ , finite or infinite, and suppose  $C_1, \dots, C_n: X \rightarrow \mathbb{R}$  is some collection of functions. Then we may seek a probability measure  $p$  that maximizes entropy subject to the constraints that the mean value of each observable  $C_i$  is a given real number  $\overline{C}_i$ :

$$\sum_{x \in X} p(x) C_i(x) = \overline{C}_i.$$

As nicely discussed by Jaynes [10, 11], the solution, if it exists, is the so-called **Gibbs ensemble**:

$$p(x) = \frac{1}{Z} e^{-(s_1 C_1(x) + \dots + s_n C_n(x))}$$

for some numbers  $s_i \in \mathbb{R}$  depending on the desired mean values  $\overline{C}_i$ . Here the normalizing factor  $Z$  is called the **partition function**:

$$Z = \sum_{x \in X} e^{-(s_1 C_1(x) + \dots + s_n C_n(x))}.$$

In thermodynamics,  $X$  represents the set of **microstates** of some physical system. A probability measure on  $X$  is also known as an **ensemble**. Each function  $C_i: X \rightarrow \mathbb{R}$  is called an **observable**, and the corresponding quantity  $s_i$  is called the **conjugate variable** of that observable. For example, the conjugate of the energy  $E$  is the inverse of temperature  $T$ , in units where Boltzmann's constant equals 1. The conjugate of the volume  $V$  — of a piston full of gas, for example — is the pressure  $P$  divided by the temperature. And in a gas containing molecules of various types, the conjugate of the number  $N_i$  of molecules of the  $i$ th type is minus the ‘chemical potential’  $\mu_i$ , again divided by temperature. For easy reference, we list these observables and their conjugate variables below.



### THERMODYNAMICS

Observable	Conjugate Variable
energy: $E$	$\frac{1}{T}$
volume: $V$	$\frac{P}{T}$
number: $N_i$	$-\frac{\mu_i}{T}$

Now let us return to the case where  $X = \text{dom}(U)$ . Recalling that programs are bit strings, one important observable for programs is the length:

$$|\cdot|: X \rightarrow \mathbb{N}.$$

We have already seen the measure

$$p(x) = \frac{1}{Z} e^{-\gamma|x|}.$$

Now its significance should be clear! This is the probability measure on programs that maximizes entropy subject to the constraint that the mean length is some constant  $\ell$ :

$$\sum_{x \in X} p(x) |x| = \ell.$$

So,  $\gamma$  is the conjugate variable to program length.

There are, however, other important observables that can be defined for programs, and each of these has a conjugate quantity. To make the analogy to thermodynamics as vivid as possible, let us arbitrarily choose two more observables and treat them as analogues of energy and the number of some type of molecule. Two of the most obvious observables are ‘output’ and ‘runtime’. Since Levin’s computable complexity measure [14] uses the logarithm of runtime as a kind of ‘cutoff’ reminiscent of an energy cutoff in renormalization, we shall arbitrarily choose the log of the runtime to be analogous to the energy, and denote it as

$$E: X \rightarrow [0, \infty)$$

Following the chart above, we use  $1/T$  to stand for the variable conjugate to  $E$ . We arbitrarily treat the output of a program as analogous to the number of a certain kind of molecule, and denote it as

$$N: X \rightarrow \mathbb{N}.$$

We use  $-\mu/T$  to stand for the conjugate variable of  $N$ . Finally, as already hinted, we denote program length as

$$V: X \rightarrow \mathbb{N}$$

so that in terms of our earlier notation,  $V(x) = |x|$ . We use  $P/T$  to stand for the variable conjugate to  $V$ .

### ALGORITHMS

Observable	Conjugate Variable
log runtime: $E$	$\frac{1}{T}$
length: $V$	$\frac{P}{T}$
output: $N$	$-\frac{\mu}{T}$

Before proceeding, we wish to emphasize that the analogies here were chosen somewhat arbitrarily. They are merely meant to illustrate the application of thermodynamics to the study of algorithms. There may or may not be a specific ‘best’ mapping between observables for programs and observables for a container of gas! Indeed, Tadaki [28] has explored another analogy, where length rather than log run time is treated as the analogue of energy. There is nothing wrong with this. However, he did not introduce enough other observables to see the whole structure of thermodynamics, as developed in Sections 4.1-4.2 below.

Having made our choice of observables, we define the partition function by

$$Z = \sum_{x \in X} e^{-\frac{1}{T}(E(x) + PV(x) - \mu N(x))}.$$

When this sum converges, we can define a probability measure on  $X$ , the Gibbs ensemble, by

$$p(x) = \frac{1}{Z} e^{-\frac{1}{T}(E(x) + PV(x) - \mu N(x))}.$$

Both the partition function and the probability measure are functions of  $T, P$  and  $\mu$ . From these we can compute the mean values of the observables to which these variables are conjugate:

$$\overline{E} = \sum_{x \in X} p(x) E(x)$$

$$\overline{V} = \sum_{x \in X} p(x) V(x)$$

$$\overline{N} = \sum_{x \in X} p(x) N(x)$$

In certain ranges, the map  $(T, P, \mu) \mapsto (\overline{E}, \overline{V}, \overline{N})$  will be invertible. This allows us to alternatively think of  $Z$  and  $p$  as functions of  $\overline{E}, \overline{V}$ , and  $\overline{N}$ . In this situation it is typical to abuse language by omitting the overlines which denote ‘mean value’.

## 4.1 Elementary Relations

The entropy  $S$  of the Gibbs ensemble is given by

$$S = - \sum_{x \in X} p(x) \ln p(x).$$

We may think of this as a function of  $T, P$  and  $\mu$ , or alternatively — as explained above — as functions of the mean values  $E, V$ , and  $N$ . Then simple calculations, familiar from statistical mechanics [22], show that

$$\left. \frac{\partial S}{\partial E} \right|_{V, N} = \frac{1}{T} \quad (3)$$

$$\left. \frac{\partial S}{\partial V} \right|_{E, N} = \frac{P}{T} \quad (4)$$

$$\left. \frac{\partial S}{\partial N} \right|_{E, V} = -\frac{\mu}{T}. \quad (5)$$

We may summarize all these by writing

$$dS = \frac{1}{T} dE + \frac{P}{T} dV - \frac{\mu}{T} dN$$

or equivalently

$$dE = T dS - P dV + \mu dN. \quad (6)$$

Starting from the latter equation we see:

$$\left. \frac{\partial E}{\partial S} \right|_{V, N} = T \quad (7)$$

$$\left. \frac{\partial E}{\partial V} \right|_{S, N} = -P \quad (8)$$

$$\left. \frac{\partial E}{\partial N} \right|_{S, V} = \mu. \quad (9)$$

With these definitions, we can start to get a feel for what the conjugate variables are measuring. To build intuition, it is useful to think of the entropy  $S$  as roughly the logarithm of the number of programs whose log runtimes, length and output lie in small ranges  $E \pm \Delta E$ ,  $V \pm \Delta V$  and  $N \pm \Delta N$ . This is at best approximately true, but in ordinary thermodynamics this approximation is commonly employed and yields spectacularly good results. That is why in thermodynamics people often say the entropy is the logarithm of the number of microstates for which the observables  $E, V$  and  $N$  lie within a small range of their specified values [22].

If you allow programs to run longer, more of them will halt and give an answer. The **algorithmic temperature**,  $T$ , is roughly the number of times

you have to double the runtime in order to double the number of ways to satisfy the constraints on length and output.

The **algorithmic pressure**,  $P$ , measures the tradeoff between runtime and length [4]: if you want to keep the number of ways to satisfy the constraints constant, then the freedom gained by having longer runtimes has to be counterbalanced by shortening the programs. This is analogous to the pressure of gas in a piston: if you want to keep the number of microstates of the gas constant, then the freedom gained by increasing its energy has to be counterbalanced by decreasing its volume.

Finally, the **algorithmic potential** describes the relation between log runtime and output: it is a quantitative measure of the principle that most large outputs must be produced by long programs.

## 4.2 Thermodynamic Cycles

One of the first applications of thermodynamics was to the analysis of heat engines. The underlying mathematics applies equally well to algorithmic thermodynamics. Suppose  $C$  is a loop in  $(T, P, \mu)$  space. Assume we are in a region that can also be coordinatized by the variables  $E, V, N$ . Then the change in **algorithmic heat** around the loop  $C$  is defined to be

$$\Delta Q = \oint_C T dS.$$

Suppose the loop  $C$  bounds a surface  $\Sigma$ . Then Stokes' theorem implies that

$$\Delta Q = \oint_C T dS = \int_{\Sigma} dT dS.$$

However, Equation (6) implies that

$$dT dS = d(T dS) = d(dE + P dV - \mu dN) = +dP dV - d\mu dN$$

since  $d^2 = 0$ . So, we have

$$\Delta Q = \int_{\Sigma} (dP dV - d\mu dN)$$

or using Stokes' theorem again

$$\Delta Q = \int_C (P dV - \mu dN). \quad (10)$$

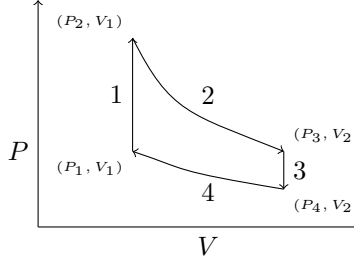
In ordinary thermodynamics,  $N$  is constant for a heat engine using gas in a sealed piston. In this situation we have

$$\Delta Q = \int_C P dV.$$

This equation says that the change in heat of the gas equals the work done on the gas — or equivalently, minus the work done *by* the gas. So, in algorithmic

thermodynamics, let us define  $\int_C PdV$  to be the **algorithmic work** done on our ensemble of programs as we carry it around the loop  $C$ . Beware: this concept is unrelated to ‘computational work’, meaning the amount of computation done by a program as it runs.

To see an example of a cycle in algorithmic thermodynamics, consider the analogue of the heat engine patented by Stoddard in 1919 [25]. Here we fix  $N$  to a constant value and consider the following loop in the  $PV$  plane:



We start with an ensemble with algorithmic pressure  $P_1$  and mean length  $V_1$ . We then trace out a loop built from four parts:

1. *Isometric*. We increase the pressure from  $P_1$  to  $P_2$  while keeping the mean length constant. No algorithmic work is done on the ensemble of programs during this step.
2. *Isentropic*. We increase the length from  $V_1$  to  $V_2$  while keeping the number of halting programs constant. High pressure means that we’re operating in a range of runtimes where if we increase the length a little bit, many more programs halt. In order to keep the number of halting programs constant, we need to shorten the runtime significantly. As we gradually increase the length and lower the runtime, the pressure drops to  $P_3$ . The total difference in log runtime is the algorithmic work done on the ensemble during this step.
3. *Isometric*. Now we decrease the pressure from  $P_3$  to  $P_4$  while keeping the length constant. No algorithmic work is done during this step.
4. *Isentropic*. Finally, we decrease the length from  $V_2$  back to  $V_1$  while keeping the number of halting programs constant. Since we’re at low pressure, we need only increase the runtime a little. As we gradually decrease the length and increase the runtime, the pressure rises slightly back to  $P_1$ . The total increase in log runtime is minus the algorithmic work done on the ensemble of programs during this step.

The total algorithmic work done on the ensemble per cycle is the difference in log runtimes between steps 2 and 4.

### 4.3 Further Relations

From the elementary thermodynamic relations in Section 4.1, we can derive various others. For example, the so-called ‘Maxwell relations’ are obtained by computing the second derivatives of thermodynamic quantities in two different orders and then applying the basic derivative relations, Equations (7-9). While trivial to prove, these relations say some things about algorithmic thermodynamics which may not seem intuitively obvious.

We give just one example here. Since mixed partials commute, we have:

$$\left. \frac{\partial^2 E}{\partial V \partial S} \right|_N = \left. \frac{\partial^2 E}{\partial S \partial V} \right|_N.$$

Using Equation (7), the left side can be computed as follows:

$$\left. \frac{\partial^2 E}{\partial V \partial S} \right|_N = \left. \frac{\partial}{\partial V} \right|_{S,N} \left. \frac{\partial E}{\partial S} \right|_{V,N} = \left. \frac{\partial T}{\partial V} \right|_{S,N}$$

Similarly, we can compute the right side with the help of Equation (8):

$$\left. \frac{\partial^2 E}{\partial S \partial V} \right|_N = \left. \frac{\partial}{\partial S} \right|_{V,N} \left. \frac{\partial E}{\partial V} \right|_{S,N} = - \left. \frac{\partial P}{\partial S} \right|_{V,N}.$$

As a result, we obtain:

$$\left. \frac{\partial T}{\partial V} \right|_{S,N} = - \left. \frac{\partial P}{\partial S} \right|_{V,N}.$$

We can also derive interesting relations involving derivatives of the partition function. These become more manageable if we rewrite the partition function in terms of the conjugate variables of the observables  $E$ ,  $V$ , and  $N$ :

$$\beta = \frac{1}{T}, \quad \gamma = \frac{P}{T}, \quad \delta = -\frac{\mu}{T}. \quad (11)$$

Then we have

$$Z = \sum_{x \in X} e^{-\beta E(x) - \gamma V(x) - \delta N(x)}$$

Simple calculations, standard in statistical mechanics [22], then allow us to compute the mean values of observables as derivatives of the logarithm of  $Z$  with respect to their conjugate variables. Here let us revert to using overlines to denote mean values:

$$\overline{E} = \sum_{x \in X} p(x) E(x) = - \frac{\partial}{\partial \beta} \ln Z$$

$$\overline{V} = \sum_{x \in X} p(x) V(x) = - \frac{\partial}{\partial \gamma} \ln Z$$

$$\overline{N} = \sum_{x \in X} p(x) N(x) = - \frac{\partial}{\partial \delta} \ln Z$$

We can go further and compute the variance of these observables using second derivatives:

$$(\Delta E)^2 = \sum_{x \in X} p(x)(E(x)^2 - \overline{E}^2) = \frac{\partial^2}{\partial^2 \beta} \ln Z$$

and similarly for  $V$  and  $N$ . Higher moments of  $E, V$  and  $N$  can be computed by taking higher derivatives of  $\ln Z$ .

#### 4.4 Convergence

So far we have postponed the crucial question of convergence: for which values of  $T, P$  and  $\mu$  does the partition function  $Z$  converge? For this it is most convenient to treat  $Z$  as a function of the variables  $\beta, \gamma$  and  $\delta$  introduced in Equation (11). For which values of  $\beta, \gamma$  and  $\delta$  does the partition function converge?

First, when  $\beta = \gamma = \delta = 0$ , the contribution of each program is 1. Since there are infinitely many halting programs,  $Z(0, 0, 0)$  does not converge.

Second, when  $\beta = 0, \gamma = \ln 2$ , and  $\delta = 0$ , the partition function converges to Chaitin's number

$$\Omega = \sum_{x \in X} 2^{-V(x)}.$$

To see that the partition function converges in this case, consider this mapping of strings to segments of the unit interval:

empty							
0				1			
00		01		10		11	
000	001	010	011	100	101	110	111
⋮							

Each segment consists of all the real numbers whose binary expansion begins with that string; for example, the set of real numbers whose binary expansion begins 0.101... is  $[0.101, 0.110)$  and has measure  $2^{-|101|} = 2^{-3} = 1/8$ . Since the set of halting programs for our universal machine is prefix-free, we never count any segment more than once, so the sum of all the segments corresponding to halting programs is at most 1.

Third, Tadaki has shown [27] that the expression

$$\sum_{x \in X} e^{-\gamma V(x)}$$

converges for  $\gamma \geq \ln 2$  but diverges for  $\gamma < \ln 2$ . It follows that  $Z(\beta, \gamma, \delta)$  converges whenever  $\gamma \geq \ln 2$  and  $\beta, \delta \geq 0$ .

Fourth, when  $\beta > 0$  and  $\gamma = \delta = 0$ , convergence depends on the machine. There are machines where infinitely many programs halt immediately. For these,  $Z(\beta, 0, 0)$  does not converge. However, there are also machines where program

$x$  takes at least  $V(x)$  steps to halt; for these machines  $Z(\beta, 0, 0)$  will converge when  $\beta \geq \ln 2$ . Other machines take much longer to run. For these,  $Z(\beta, 0, 0)$  will converge for even smaller values of  $\beta$ .

Fifth and finally, when  $\beta = \gamma = 0$  and  $\delta > 0$ ,  $Z(\beta, \gamma, \delta)$  fails to converge, since there are infinitely many programs that halt and output 0.

## 4.5 Computability

Even when the partition function  $Z$  converges, it may not be computable. The theory of computable real numbers was independently introduced by Church, Post, and Turing, and later blossomed into the field of computable analysis [21]. We will only need the basic definition: a real number  $a$  is **computable** if there is a recursive function that maps any natural number  $n > 0$  to an integer  $f(n)$  such that

$$\frac{f(n)}{n} \leq a \leq \frac{f(n) + 1}{n}.$$

In other words, for any  $n > 0$ , we can compute a rational number that approximates  $a$  with an error of at most  $1/n$ . This definition can be formulated in various other equivalent ways: for example, the computability of binary digits.

Chaitin [6] proved that the number

$$\Omega = Z(0, \ln 2, 0)$$

is uncomputable. In fact, he showed that for any universal machine, the values of all but finitely many bits of  $\Omega$  are not only uncomputable, but random: knowing the value of some of them tells you nothing about the rest. They're independent, like separate flips of a fair coin.

More generally, for any computable number  $\gamma \geq \ln 2$ ,  $Z(0, \gamma, 0)$  is ‘partially random’ in the sense of Tadaki [3, 27]. This deserves a word of explanation. A fixed formal system with finitely many axioms can only prove finitely many bits of  $Z(0, \gamma, 0)$  have the values they do; after that, one has to add more axioms or rules to the system to make any progress. The number  $\Omega$  is completely random in the following sense: for each bit of axiom or rule one adds, one can prove at most one more bit of its binary expansion has the value it does. So, the most efficient way to prove the values of these bits is simply to add them as axioms! But for  $Z(0, \gamma, 0)$  with  $\gamma > \ln 2$ , the ratio of bits of axiom per bits of sequence is less than 1. In fact, Tadaki showed that for any computable  $\gamma \geq \ln 2$ , the ratio can be reduced to exactly  $(\ln 2)/\gamma$ .

On the other hand,  $Z(\beta, \gamma, \delta)$  is computable for all computable real numbers  $\beta > 0$ ,  $\gamma \geq \ln 2$  and  $\delta \geq 0$ . The reason is that  $\beta > 0$  exponentially suppresses the contribution of machines with long runtimes, eliminating the problem posed by the undecidability of the halting problem. The fundamental insight here is due to Levin [14]. His idea was to ‘dovetail’ all programs: on turn  $n$ , run each of the first  $n$  programs a single step and look to see which ones have halted. As they halt, add their contribution to the running estimate of  $Z$ . For any  $k \geq 0$  and turn  $t \geq 0$ , let  $k_t$  be the location of the first zero bit after position  $k$  in the



estimation of  $Z$ . Then because  $-\beta E(x)$  is a monotonically decreasing function of the runtime and decreases faster than  $k_t$ , there will be a time step where the total contribution of all the programs that have not halted yet is less than  $2^{-k_t}$ .

## 5 Conclusions

There are many further directions to explore. Here we mention just three. First, as already mentioned, the ‘Kolmogorov complexity’ [12] of a number  $n$  is the number of bits in the shortest program that produces  $n$  as output. However, a very short program that runs for a million years before giving an answer is not very practical. To address this problem, the **Levin complexity** [15] of  $n$  is defined using the program’s length plus the logarithm of its runtime, again minimized over all programs that produce  $n$  as output. Unlike the Kolmogorov complexity, the Levin complexity is computable. But like the Kolmogorov complexity, the Levin complexity can be seen as a *relative entropy*—at least, up to some error bounded by a constant. The only difference is that now we compute this entropy relative to a different probability measure: instead of using the Gibbs distribution at infinite algorithmic temperature, we drop the temperature to  $\ln 2$ . Indeed, the Kolmogorov and Levin complexities are just two examples from a continuum of options. By adjusting the algorithmic pressure and temperature, we get complexities involving other linear combinations of length and log runtime. The same formalism works for complexities involving other observables: for example, the maximum amount of memory the program uses while running.

Second, instead of considering Turing machines that output a single natural number, we can consider machines that output a finite list of natural numbers  $(N_1, \dots, N_j)$ ; we can treat these as populations of different “chemical species” and define algorithmic potentials for each of them. Processes analogous to chemical reactions are paths through this space that preserve certain invariants of the lists. With chemical reactions we can consider things like internal combustion cycles.

Finally, in ordinary thermodynamics the partition function  $Z$  is simply a number after we fix values of the conjugate variables. The same is true in algorithmic thermodynamics. However, in algorithmic thermodynamics, it is natural to express this number in binary and inquire about the algorithmic entropy of the first  $n$  bits. For example, we have seen that for suitable values of temperature, pressure and chemical potential,  $Z$  is Chaitin’s number  $\Omega$ . For each universal machine there exists a constant  $c$  such that the first  $n$  bits of the number  $\Omega$  have at least  $n - c$  bits of algorithmic entropy with respect to that machine. Tadaki [27] generalized this computation to other cases.

So, *in algorithmic thermodynamics, the partition function itself has nontrivial entropy*. Tadaki has shown that the same is true for algorithmic pressure (which in his analogy he calls ‘temperature’). This reflects the self-referential nature of computation. It would be worthwhile to understand this more deeply.

## Acknowledgements

We thank Leonid Levin and the denizens of the  $n$ -Category Café for useful comments. MS thanks Cristian Calude for many discussions of algorithmic information theory. JB thanks Bruce Smith for discussions on relative entropy. He also thanks Mark Smith for conversations on physics and information theory, as well as for giving him a copy of Reif's *Fundamentals of Statistical and Thermal Physics*.

## References

- [1] C. H. Bennett, P. Gacs, M. Li, M. B. Vítányi and W. H. Zurek, Information distance, *IEEE Trans. Inform. Theor.* **44** (1998), 1407–1423.
- [2] C. S. Calude, *Information and Randomness: An Algorithmic Perspective*, Springer, Berlin, 2002.
- [3] C. S. Calude, L. Staiger, S. A. Terwijn, On partial randomness, *Ann. Appl. Pure Logic* **138** (2006) 20–30. Also available at <http://www.cs.auckland.ac.nz/CDMTCS/researchreports/239cris.pdf>.
- [4] C. S. Calude and M. A. Stay, Most Programs Stop Quickly or Never Halt, *Adv. Appl. Math.* **40** (3), 295–308. Also available as [arXiv:cs/0610153](http://arxiv.org/abs/cs/0610153).
- [5] C. S. Calude and M. A. Stay, Natural halting probabilities, partial randomness, and zeta functions, *Inform. and Comput.*, **204** (2006), 1718–1739.
- [6] G. Chaitin, A theory of program size formally identical to information theory, *Journal of the ACM* **22** (1975), 329–340. Also available at <http://www.cs.auckland.ac.nz/~chaitin/acm75.pdf>.
- [7] G. Chaitin, Algorithmic entropy of sets, *Comput. Math. Appl.* **2** (1976), 233–245. Also available at <http://www.cs.auckland.ac.nz/CDMTCS/chaitin/sets.ps>.
- [8] C. P. Roberts, *The Bayesian Choice: From Decision-Theoretic Foundations to Computational Implementation*, Springer, Berlin, 2001.
- [9] E. Fredkin and T. Toffoli, Conservative logic, *Intl. J. Theor. Phys.* **21** (1982), 219–253. Also available at <http://strangepaths.com/wp-content/uploads/2007/11/conservativelogic.pdf>.
- [10] E. T. Jaynes, Information theory and statistical mechanics, *Phys. Rev.* **106** (1957), 620–630. Also available at <http://bayes.wustl.edu/etj/articles/theory.1.pdf>.
- [11] E. T. Jaynes, *Probability Theory: The Logic of Science*, Cambridge U. Press, Cambridge, 2003. Draft available at <http://omega.albany.edu:8008/JaynesBook.html>.

- [12] A. N. Kolmogorov, Three approaches to the definition of the quantity of information, *Probl. Inf. Transm.* **1** (1965), 3–11.
- [13] S. Kullback and R. A. Leibler, On information and sufficiency, *Ann. Math. Stat.* **22** (1951), 79–86.
- [14] L. A. Levin, Universal sequential search problems, *Probl. Inf. Transm.* **9** (1973), 265–266.
- [15] L. A. Levin, Laws of information conservation (non-growth) and aspects of the foundation of probability theory. *Probl. Inf. Transm.* **10** (1974), 206–210.
- [16] L. A. Levin and A. K. Zvonkin, The complexity of finite objects and the development of the concepts of information and randomness by means of the theory of algorithms, *Russian Mathematics Surveys* **256** (1970), 83–124  
Also available at <http://www.cs.bu.edu/fac/lnd/dvi/ZL-e.pdf>.
- [17] M. Li and P. Vitányi, *An Introduction to Kolmogorov Complexity Theory and its Applications*, Springer, Berlin, 2008.
- [18] Y. Manin, *Renormalization and computation I: motivation and background*. Available as arxiv:0904.4921.
- [19] Y. Manin, *Renormalization and computation II: Time Cut-off and the Halting Problem*. Available as arxiv:0908.3430.
- [20] Y. Manin, M. Marcolli, *Error-correcting codes and phase transitions*. Available as arxiv:0910.5135.
- [21] M. B. Pour-El and J. I. Richards, *Computability in Analysis and Physics*, Springer, Berlin, 1989. Also available at <http://projecteuclid.org/euclid.pl/1235422916>.
- [22] F. Reif, *Fundamentals of Statistical and Thermal Physics*, McGraw-Hill, New York, 1965.
- [23] A. Rényi, On measures of information and entropy, *Proceedings of the 4th Berkeley Symposium on Mathematics, Statistics and Probability*, 1960, pp. 547–561. Also available at [http://digitalassets.lib.berkeley.edu/math/ucb/text/math\\_s4\\_v1\\_article-27.pdf](http://digitalassets.lib.berkeley.edu/math/ucb/text/math_s4_v1_article-27.pdf).
- [24] R. J. Solomonoff, A formal theory of inductive inference, part I, *Inform. Control* **7** (1964), 1–22. Also available at <http://world.std.com/~rjs/1964pt1.pdf>.
- [25] E. J. Stoddard, Apparatus for obtaining power from compressed air, US Patent 1,926,463. Available at <http://www.google.com/patents?id=zLRFAAAAEBAJ>.

- [26] L. Szilard, On the decrease of entropy in a thermodynamic system by the intervention of intelligent beings, *Zeit. Phys.* **53** (1929) 840–856. English translation in H. S. Leff and A. F. Rex (eds.) *Maxwell's Demon. Entropy, Information, Computing*, Adam Hilger, Bristol, 1990.
- [27] K. Tadaki, A generalization of Chaitin's halting probability  $\Omega$  and halting self-similar sets, *Hokkaido Math. J.* **31** (2002), 219–253. Also available as arXiv:nlin.CD/0212001.
- [28] K. Tadaki, A statistical mechanical interpretation of algorithmic information theory. Available as arXiv:0801.4194.
- [29] K. Tadaki, A statistical mechanical interpretation of algorithmic information theory III: Composite systems and fixed points. *Proceedings of the 2009 IEEE Information Theory Workshop, Taormina, Sicily, Italy*, to appear. Also available as arXiv:0904.0973.