**Review article**

# Quantifying division of labor: borrowing tools from sociology, sociobiology, information theory, landscape ecology, and biogeography

**R. Gorelick[1,2] and S.M. Bertram[1]**

[1] *Department of Biology, Carleton University, Ottawa, Ontario K1S 5B6, Canada <u>and</u> School of Life Sciences, Arizona State University, Tempe, AZ 85287–4501, U.S.A.; e-mail: Root_Gorelick@carleton.ca, Susan_Bertram@carleton.ca*
[2] *U.S. Environmental Protection Agency, National Center for Environmental Assessment, Washington DC 20460, U.S.A.*

**Abstract.** How do we quantify division of labor? We review several fields (sociology, landscape ecology, statistics, information theory, and biogeography) that have been cognizant of these questions and been somewhat successful at answering them. We review fourteen indices for quantifying division of labor, *sensu lato*, which can be categorized into four families: Shannon's index/entropy, Simpson's index, geometric mean, and standard/absolute deviation (including coefficients of variation). We argue that those indices using matrix inputs will simultaneously quantify the interplay between all individuals and all tasks and will thus best capture the essence of division of labor.

*Keywords:* Sociality, mutual information, task specialization, information theory.

## Introduction

Division of labor is one of the primary attributes of sociality; it is thought to be ubiquitous across social groups, and of paramount importance in explaining the success of eusocial organisms (Wilson, 1971). Quantifying division of labor is, therefore, of central importance to sociobiology (Beshers and Fewell, 1996). Quantifying division of labor is also important to fields outside of sociobiology, fields that include economics, sociology, landscape ecology, statistics, information theory, biogeography, and conservation. These fields have all independently constructed indices that quantify their notion of division of labor. Here we review and compare the statistics that have arisen from these fields. Our interdisciplinary approach allows us to re-conceptualize not

only individuals and tasks as other constructs, but it also allows us to re-conceptualize the data linking individuals to tasks to be something other than how much time an individual spends performing each given task. We explicate each field's equivalences to division of labor in greater detail below. By using the statistical tools that have already been developed by these fields, scientists may be able to consistently analyze the complex data that arises when studying division of labor in a way that will allow it to be used by researchers across fields.

Sociologists and economists take the identical view as sociobiologists as to what constitutes division of labor – where different individuals perform different tasks, and some individuals specialize on certain tasks (Michener, 1974). Most other fields, however, require us to re-conceptualize how we think about division of labor. Information theorists, for example, use division of labor in a manner identical to that of sociobiologists, but they refer to it as mutual information. Mutual information describes how different bits of a message are transmitted in both directions between two parties (Shannon, 1948). As we describe in the next paragraph, from the division of labor perspective, each sent bit is equivalent to an individual and each received bit is equivalent to a task.

Shannon had originally constructed mutual entropy to simultaneously measure the amount of information transmitted in both directions between two parties, hence the moniker mutual information. This is exactly what we want in measuring division of labor. Here, the first party consists of the ensemble of individuals in the population and the second party consists of the ensemble of tasks. Assume that division of labor is high. Then, for each individual, we should be able to predict which task it is performing. High division of labor means that information about the individuals gets transmitted to infor-

**Table 1.** Division of labor formulae

---

Shannon/Entropy

1a. $H = -\sum\limits_{i=1}^{n} p_i \cdot \ln(p_i)$

1b. $I = \sum\limits_{\substack{i=1 \\ j=1}}^{m,n} p_{ij} \cdot \ln\left(\frac{p_{ij}}{p_i \cdot p_j}\right)$

1c. $I/H$

Simpson

2a. $S = \sum\limits_{i=1}^{n} p_i^2$

2b. $1 - S = 1 - \sum\limits_{i=1}^{n} p_i^2$

2c. $\frac{1}{S} = \left[\sum\limits_{i=1}^{n} p_i^2\right]^{-1}$

2d. $-\ln(S) = -\ln\left(\sum\limits_{i=1}^{n} p_i \cdot p_i\right)$

2e. $\varsigma = -\ln\left(\sum\limits_{\substack{i=1 \\ j=1}}^{n,m} \frac{p_{ij}}{p_i} \cdot \frac{p_{ij}}{p_j}\right)$

2f. $\varsigma/S$

Geometric mean

3. $G = \exp\left[\frac{1}{n}\sum\limits_{i=1}^{n} \ln(p_i)\right]$

Standard/absolute deviations

4a. $D = \dfrac{\left(\sum\limits_{i=1}^{n} p_i\right)^2 - \left(\sum\limits_{i=1}^{n} p_i^2\right)}{\left(\sum\limits_{i=1}^{n} p_i\right)^2} = 1 - \dfrac{\left(\sum\limits_{i=1}^{n} p_i^2\right)}{\left(\sum\limits_{i=1}^{n} p_i\right)^2}$

4b. $D' = \frac{D}{D_{\max}} = \frac{D}{\left(1 - \frac{1}{n}\right)}$

4c. $1 - \dfrac{\left(SD(p)/\bar{p}\right)}{\left(SD(p)/\bar{p}\right)_{\max}} = 1 - \dfrac{SD(p)/\bar{p}}{\sqrt{n(n-1)}} = 1 - \sqrt{\frac{1}{n(n-1)}\left[\sum\limits_{i=1}^{n}(p_i - \bar{p})^2\right]}$ where $\bar{p} = \sum\limits_{i=1}^{n} p_i$

4d. $1 - \dfrac{\left(AD(p)/\bar{p}\right)}{\left(AD(p)/\bar{p}\right)_{\max}} = 1 - \dfrac{1}{2 \cdot \left(1 - \frac{1}{n}\right)} \cdot \dfrac{AD(p)}{\bar{p}} = 1 - \dfrac{1}{2 \cdot \left(1 - \frac{1}{n}\right)} \cdot \dfrac{\sum\limits_{i=1}^{n}|p_i - \bar{p}|}{\bar{p}}$

---

mation about tasks. Likewise, with high division of labor, because each task is performed by so few individuals, information about which task is being performed gets transmitted into information about which individual is performing that task. This is bidirectional information transmission, exactly as envisioned by Shannon, with nothing but new interpretations of the variables. Conversely, now assume that division of labor is low. Then knowing which task was selected from the ensemble transmits no information about which individual is most likely performed it <u>or</u> knowing which individual you encounter transmits no information about which task it is performing (N.B. note the conjunction "or" and not "and" here). Information transmission is between individuals and tasks, and vice versa.

Biogeography also quantifies division of labor regularly in their measure of biodiversity. Here, biodiversity takes into account the numbers of each species in each patch (i.e., individuals and tasks, respectively; Gorelick, 2006). From a conservation biology perspective, common species that are endemic to single locales (high division of labor) can be as much of a conservation concern as rare cosmopolitan species (low division of labor). Rare cosmopolitan species are analogous to rare task special-

ists, such as undertakers in an ant colony (Julian and Cahan, 1999). For landscape ecologists, the equivalent to division of labor is redundancy or $\beta$ diversity, which quantifies the interdependence between two measures of geographic patches (Ernoult et al., 2003). For example, is land-use type driven by edaphic conditions (high division of labor) or is land-use type independent of edaphic conditions (low division of labor)?

We review four classes of division of labor statistics (Table 1). First, we discuss three statistics (1a–c) that are based on Shannon's index, then six statistics (2a–f) that are based on Simpson's index, one statistic (3) based on the geometric mean, and finally four interrelated statistics that are based on the standard/absolute deviation (including coefficients of variation) (indices 4a–d).

These fourteen indices can be categorized in an alternative, orthogonal direction: whether their inputs are vectors or matrices. Because we envision division of labor as reciprocal communication between individuals and tasks (alternatively, think of employers and employees), the indices with matrix inputs will prove to be the most valuable.

## Notation

We use $p_i$ to represent the probability of event $i$ and $p_{ij}$ to represent the joint probability of events $i$ and $j$. Being probabilities, each individual value of $p_i$ and $p_{ij}$ is greater than or equal to zero and $\sum_i p_i = \sum_{i,j} p_{ij} = 1$. We use $i$ to index the $n$ individuals and $j$ to index the $m$ tasks. Probabilities can be interpreted as proportions or frequency of occurrence of events, e.g. $p_{ij}$ is the proportion of total time that individual $i$ expends on task $j$. Inputs to each of the division of labor statistics herein will be either the vector $[p_i] = (p_1, p_2, \ldots, p_i, \ldots p_n)$ or the matrix

$$[p_{ij}] = \begin{pmatrix} p_{11} & p_{12} & \cdots & p_{1n} \\ p_{21} & p_{21} & \cdots & p_{2n} \\ \vdots & \vdots & & \vdots \\ p_{m1} & p_{m2} & \cdots & p_{mn} \end{pmatrix}. \text{ If data are not originally}$$

given in terms of probabilities, then take each data point and divide it by $\sum_i^n p_i$ or $\sum_{i,j}^{n,m} p_{ij}$, depending on whether the input is a vector or matrix, respectively.

## Indices based on Shannon's index

Communication between two parties was the basis for Shannon's pioneering work on mutual entropy (1948). Biologists refer to marginal entropy (not mutual entropy) as the 'Shannon index' and use this as a measure of division of labor.

Index 1a is marginal entropy $H = -\sum_{i=1}^n p_i \cdot \ln(p_i)$, which was first introduced by Boltzmann (1872) in statistical physics for quantifying the energy in an ensemble of $n$ particles. Marginal entropy is probably the most commonly used measure of division of labor. Marginal entropy has the advantage of being a concave function (Lande, 1996). Concavity means that if we take two separate populations and aggregate them, the division of labor for the aggregated metapopulation should be at least as great as the average division of labor of the two individual populations.

Kolmes (1985) used marginal entropy to quantify how specialized individuals are on tasks. Kolmes computed marginal entropy for each individual and then computed the arithmetic mean for all individuals in a population. While this approach quantifies division of labor, a problem unfortunately arises when all individuals are specialists ($p_i = 1$ for all individuals $i$), but all individuals specialize on the same task. In this case, the mean marginal entropy would be maximal, but heuristically the division of labor statistic should yield a minimum value because all individuals are performing the same task.

Index 1b is mutual entropy $I = \sum_{\substack{i=1 \\ j=1}}^{m,n} p_{ij} \cdot \ln\left(\frac{p_{ij}}{p_i \cdot p_j}\right)$. Shannon (1948) showed great insight into quantifying division of labor when he introduced mutual entropy ($I$). Mutual entropy shows how we can take a matrix of data that describes the proportion of time each individual spends performing each task and distill it into a single number. For Shannon's communications theory applications, this number represented the amount of information transferred between two parties through a noisy communications channel. Both parties act simultaneously as senders and receivers, much as people do during a typical phone conversation listening and responding to each other. From the sociobiology perspective, the role of one party to the phone conversation is filled by all the individuals in the population and the role of the other party is filled by all the tasks that they perform. Mutual entropy can therefore be used to quantify division of labor when the matrix is made up of individuals and tasks (Gorelick et al., 2004). Mutual entropy also quantifies the interdependency between rows and columns of a matrix. As such, mutual entropy has been used to quantify niche overlap (Colwell and Futuyma, 1971), stratification of data (Martínez et al., 1992), or ecological $\beta$-diversity (in the sense of how do soil conditions and land-use types influence on another; Ernoult et al., 2003).

Index 1c normalizes mutual entropy $D = \frac{I}{H}$. Mutual entropy, $I$, can be as big as the logarithm of the number or rows or columns of the data matrix. Therefore, as the number of individuals increases, so will mutual entropy. This is not good, insofar as we often want to measure whether division of labor changes as population size changes. Therefore, Gorelick et al. (2004) proposed to divide mutual entropy by the marginal entropy of the individuals, which here is equivalent to the logarithm of population size. Normalizing the division of labor statistic forces it to take on a finite maximum value of one. This statistic was first introduced in sociobiology to measure division of labor (Gorelick et al., 2004) and is now also being used in ecology to quantify $\beta$-diversity (Gorelick, 2006). Refer to Gorelick et al. (2004) and Gorelick (2006) for a set of comprehensive examples which reveal how normalized mutual entropy works to quantify division of labor.

## Indices based on Simpson's index

Index 2a is the classical Simpson's index, $S = \sum_{i=1}^n p_i^2$, i.e. the sum of squares of a vector of probabilities, were $1 = \sum_{i=1}^n p_i$ because these are probabilities. Simpson's index was first introduced to measure ecological diversity, where the data was in the form of abundances of a set of species, but could also be used as a proxy for division of labor.

Indices 2b, 2c, and 2d are minor variants on Simpson's index. Each has been used as a measure of ecological diversity and each is directly proportional to diversity. Note that diversity is inversely proportional to division of labor, which was the likely justification for each of these modifications. Therefore Indices 2b, 2c, and 2d will be inversely proportional to division of labor.

Index 2b is one minus Simpson's index, $1 - \sum_{i=1}^{n} p_i^2$. It equals the probability that two random individuals perform the same task (Lande, 1996), and is often called the Gini index after its inventor (Gini, 1912), who preceded Simpson by several decades. Similar to marginal entropy (index 1a), the Gini index is also concave.

Index 2c is the reciprocal of Simpson's index, $\left[\sum_{i=1}^{n} p_i^2\right]^{-1}$. It is not concave, but arises naturally as a special case of the Hill family, $\left[\sum_{i=1}^{n} p_i^k\right]^{\frac{1}{1-k}}$, when $k{=}2$. A transformation of mutual entropy also is part of the Hill family (Keylock, 2005).

Index 2d is the logarithm of Simpson's index, $-\ln\left(\sum_{i=1}^{n} p_i^2\right)$. Buckland et al. (2005) prefer this index for quantifying division of labor because it has an estimator whose expected value does not depend on sample size. We also prefer this index because it has the same functional form as marginal entropy (Index 1a) and hence can be readily generalized to matrix inputs.

The problem with Simpson's index and the above generalizations of it (indices 2a-d) is identical to what we discussed with indices 1a and 1b. They can only be computed for a single task or single individual at a time and not for the entire ensemble of tasks being done by all individuals.

Index 2e has matrices (rather than vectors) as input and is the matrix analogue of Simpson's index $\varsigma = -\ln\left(\sum_{\substack{i=1 \\ j=1}}^{n,m} \frac{p_{ij}}{p_i} \cdot \frac{p_{ij}}{p_j}\right)$ (Gorelick, 2006). As with the normalized mutual entropy index (1c), this quantity can be normalized by dividing by the ordinary Simpson's index. As with Shannon's index, we divide by whichever is larger: Simpson's index of $p_i$ or $p_j$ (explained in greater detail in the section on comparing indices, below).

### Index based on the geometric mean

Index 3 uses the geometric mean of the elements of the probability vector to quantify division of labor for a single individual (Buckland et al., 2005): $G = \exp\left[\frac{1}{n}\sum_{i=1}^{n} \ln(p_i)\right]$. This index also suffers from the problem that it does not simultaneously use all the information in the data matrix of individuals and tasks. It was first introduced

to quantify ecological diversity of abundances of several species.

### Interrelated indices

The remaining four interrelated division of labor statistics come from the sociology literature and were originally published several decades ago. They are approximately coefficients of variation, but where the numerator can be either standard deviation or average deviation and where the denominator can be either the mean or its square. Neither these papers (Gibbs and Martin, 1962; Labovitz and Gibbs, 1964; Gibbs and Poston, 1965; Martin and Gray, 1971) nor their immediate descendants (e.g. Rushing and Davies, 1970; Smith and Snow, 1976) are regularly cited by sociobiologists. This suggests that these indices are mostly unknown by sociobiologists. All four of these indices use the probability vector from each individual or task as the input and then compute some normalized form of the standard deviation or average deviation for each probability vector; none use matrix inputs. When these indices where originally created, it was not obvious they were related to one another. By rewriting the formulae in terms of standard deviation and average deviation ($SD$ and $AD$), the relationships are much more readily evident.

Index 4a divides the variance of the probabilities by the square of the mean probability (Gibbs and Martin, 1962): $D = \dfrac{\left(\sum_{i=1}^{n} p_i\right)^2 - \left(\sum_{i=1}^{n} p_i^2\right)}{\left(\sum_{i=1}^{n} p_i\right)^2} = 1 - \dfrac{\left(\sum_{i=1}^{n} p_i^2\right)}{\left(\sum_{i=1}^{n} p_i\right)^2}$.

Even though $D$ is bounded, it may be greater than one. Therefore, to normalize it, divide by the maximum possible value of $D$, which will depend on the length of probability vector, i.e. the number of tasks an individual performs (Labovitz and Gibbs, 1964; Rushing and Davies, 1970). It turns out that $D_{max} = 1 - \frac{1}{n}$, where $n$ is the length of the vector. Therefore, the normalized version of $D$ is given by index 4b, which is $D' = \frac{D}{\left(1-\frac{1}{n}\right)}$.

Index 4c uses the coefficient of variation (CV) to quantify division of labor. The coefficient of variation normalizes standard deviation by dividing it by the mean. Although $1 - CV$ would provide a reasonable estimate of division of labor, we have never seen this used. Instead, division of labor has been quantified using CV by further normalizing to account for length of the probability vector as $1 - \frac{CV}{CV_{max}} = 1 - \frac{CV}{\sqrt{n(n-1)}}$ (Martin and Gray, 1971; Smith and Snow, 1976).

The normalized coefficient of variation, index 4c, has also been used to measure linearity of dominance hierarchies in animal behavior, which can be considered a division of labor (Landau, 1951). A strict dominance hierarchy is a highly structured society whose division of

labor value (a.k.a. Landau linearity index) attains a maximum possible value of one. The input data is usually envisioned as a vector with one component for each individual, where the value of each component represents the number of other individuals that the focal individual is dominant over, but could equally well represent the amount of time or effort each individual spends performing a particular task. More recently, the input data has been envisioned as a square matrix, whose values represent how much individual $i$ dominates individual $j$ (de Vries, 1995). Values in the matrix can be binary or any non-negative number. However, sums of each row are computed to form a vector which is used to compute a normalized coefficient of variation, i.e. this is fundamentally a vector input.

Finally, index 4d uses absolute deviations. Instead of starting with standard deviations (i.e. Euclidean metric), some authors (Gibbs and Poston, 1965; Smith and Snow, 1976) have started with absolute deviation (i.e. taxi cab metric). They then computed the coefficient of variation as the absolute deviation (AD) divided by the mean to yield the following counterpart to the previous division of labor metric:

$$1 - \frac{\left(AD(p)_{/\bar{p}}\right)}{\left(AD(p)_{/\bar{p}}\right)_{max}} = 1 - \frac{AD(p)}{2 \cdot \left(1 - \frac{1}{n}\right) \cdot \bar{p}}.$$

### Other indices

These fourteen indices are not the only division of labor statistics in the literature. Other, more specialized indices have been developed, but these only appear to be useful for rather specific data structures. We present three examples. Gautrais et al. (2002) define division of labor for a binary state variable by how often that state changes. Duncan and Duncan (1955) define division of labor amongst women and men, but using a statistic that cannot be extended to more than two sexes, i.e. cannot be extended to division of labor if there are more than two individuals in the population. Division of labor of reproduction can be quantified as reproductive skew, which is an amalgamation of the standard deviation of reproductive output of breeding individuals weighted by the number of breeders and non-breeders (Reeve and Keller, 1996). But it is not obvious that this index can be generalized to situations with more than two tasks (e.g. more than just breeding and non-breeding). However, we have provided (Gorelick et al., 2004) an index of reproductive skew using mutual entropy to computing division of labor by letting rows of the matrix be the females in a population, columns in the matrix be males in the population, and the matrix elements be how often a given female and male mate (or a binary variable of whether they mate). Curiously, considering reproductive skew in this way – as a matrix of interactions between females and males – is commensurate with the notion of sex as a form of communication or social interaction (Roughgarden, 2004). By contrast, the fourteen indices

described above all have as inputs the frequency (probability) with which individual performs each task, without necessarily a temporal sequence of task performance and without a restriction to a pair of individuals or tasks.

### Comparing indices

The most fundamental difference between the various division of labor indices is whether they operate on one vector – i.e. one individual or one task – at a time or instead simultaneously operate on an entire matrix of data of all individuals and tasks. In economics, the matrix approach is known as a Leontief input-output model (1951). Only the mutual information and the matrix analogue of Simpson's index simultaneously use a matrix approach (Indices 1b–c and 2e–f, respectively). These matrix statistics have rarely been used by biologists, although biogeochemists have utilized this matrix approach. Bormann and Likens (1967), for example, cited Leontief's pioneering work, while Ernoult et al. (2003) used mutual information to measure the interdependencies between inputs and outputs.

Using the vector statistics (1a, 2a–d, 3, 4a–d), one can create a single statistic for all individuals or all tasks. For example, one could compute division of labor for each individual and then compute the geometric mean across all individuals (e.g. Kolmes, 1985). The same can be done for all tasks (e.g., Seeley, 1982; Beshers and Traniello, 1992). But, as Rushing (1968) and others have noted, division of labor has two components: structural/societal differentiation and individual differentiation. These two components correspond with division of labor and division of task. Unfortunately, computing these components separately does not account for interactions between individuals and tasks. For example, Kolmes (1985) approach of computing the arithmetic mean of Shannon's index across all individuals cannot distinguish between the following two populations in which each individual specializes on a single task. In the first population, half of the individuals specialize on one task and half specialize on a second task. In the second population, one individual specializes on one task and all the other individuals specialize on the second task. The arithmetic mean of Shannon's index yields maximum division of labor for both populations, yet we need to somehow account for greater division of labor in the first population in which individuals do a better job of divvying up the two tasks. The only way to quantify interdependencies between individuals (and between tasks) with vector statistics is to report multiple scalar indices. For example compute the ordered pair consisting of the arithmetic mean of Shannon's index of the individuals and the arithmetic mean of Shannon's index of the tasks. The biggest problem with this approach is that one cannot determine whether one ordered pairs is larger than another ordered pair. That is, you cannot determine whether one population has larger division of labor than another. We should

not, however, dismiss the vector statistics. They are useful so long as we are willing to abandon the criterion of a single (well-ordered) scalar statistic.

There seem to be few good reasons to favor one vector statistic over another. They are all conceptually simple, even if the formulae for some look less intuitive than others. Some are normalized for length of the input vector (i.e. for numbers of individuals or tasks), while others are not. However, they all could be readily normalized. We therefore pass no judgment over which vector statistic is most appropriate in general nor for a specific application.

To allow a more universal understanding of what is meant by division of labor, we ask that the readers re-conceptualize division of labor to mean reciprocal and cooperative *communication*, where communication is any signal or action that one individual transmits to others. This interpretation of division of labor should not be considered too unusual insofar as division of labor was defined that way by E.O. Wilson in 1971 and more recently by Costa and Fitzgerald in 1996. This conceptual leap is homologous with that made by Claude Shannon (1948) in communications theory almost six decades ago. With highly divided labor, not only do individuals divvy up tasks and thereby each individual becomes a specialist, but also the tasks are divided up amongst the individuals. If we think of tasks as employers, then high division of labor means that individuals each only work for one employer and employers specialize in who they hire.

The homologous conceptualization of division and labor with reciprocal communication can be made even more precise with elementary mathematics. Division of labor and reciprocal communication *sensu* Shannon both have matrices as inputs. For division of labor, the rows are individuals and the columns are tasks. For communication signals, the rows and columns are bits of information received and sent between a pair of individuals. High division of labor means that given an individual (row of the matrix), we can predict what task it will perform (column). High division of labor also means that given a task (column), we can predict which individual performed it (row). This is the crux of mutual entropy and its analogue for Simpson's index.

We thus strongly encourage researchers to use the matrix statistics in lieu of vector statistics insofar as only the former can capture interdependencies between all individuals and all tasks in a single scalar statistic. But, which one of these matrix statistics should a researcher use: 1b, 1c, 2e or 2f – $I = \sum_{\substack{i=1 \\ j=1}}^{m,n} p_{ij} \cdot \ln\left(\frac{p_{ij}}{p_i \cdot p_j}\right)$, $I/H$,

$\varsigma = -\ln\left(\sum_{\substack{i=1 \\ j=1}}^{n,m} \frac{p_{ij}}{p_i} \cdot \frac{p_{ij}}{p_j}\right)$, or $\varsigma/S$ ? We previously discussed

the first two of these statistics at length in Gorelick et al. (2004) and the latter two in Gorelick (2006). $I$ is not normalized for length of the probability vector, i.e. not normalized for the number of individuals or tasks. There

are three forms of $I/H$ because $H$ can be the marginal entropy of the individuals or tasks or the geometric mean of the two. Researchers should divide by marginal entropy of individuals, if there are more individuals than tasks… and vice versa if there are more tasks than individuals. Dividing by the geometric mean of the marginal entropy of individuals and tasks only makes sense if individuals and tasks are re-conceptualized to represent roughly equivalent entities. For example, if one were interested in sex ratios, females could be considered the individuals and males could be considered the tasks. Researchers could divide by the geometric mean of their marginal entropies if there was roughly a constant ratio of females to males across populations that were being compared.

The only place where we have seen a matrix version of Simpson's index used is in computing biodiversity, where Rao (1982) introduced quadratic Simpson's index $s = \sum_{\substack{i=1 \\ j=1}}^{n,m} \rho_{ij} \cdot p_i \cdot p_j$, where $\rho_{ij}$ is the taxonomic distance between species $i$ and $j$, while $p_i$ and $p_j$ are probabilities of finding species $i$ and $j$. This simply is not useful for division for labor because there is no such entity as distance between individuals and tasks.

To our knowledge, we are the first to discuss Simpson's index in the context of division of labor and to compute

$\varsigma = -\ln\left(\sum_{\substack{i=1 \\ j=1}}^{n,m} \frac{p_{ij}}{p_i} \cdot \frac{p_{ij}}{p_j}\right)$ or $\varsigma/S$ (indices 2e and 2f, respective-

ly). There are identical issues with normalizing Simpson's and Shannon's indices. There are three different possible normalizations – $\varsigma/S(p_i)$, $\varsigma/S(p_j)$, and $\frac{\varsigma}{\sqrt{S(p_i) \cdot S(p_j)}}$ – and the biology dictates the choice (Gorelick et al., 2004) in the exact same way we prescribed for normalized mutual entropy (two paragraphs above). For traditional division of labor applications, if there are more individuals than tasks, then use $\varsigma/S(p_i)$ or $I/H(p_i)$. If there are more tasks than individuals, then use $\varsigma/S(p_j)$ or $I/H(p_j)$. If the ratio of individuals to task remains constant, then use $\frac{\varsigma}{\sqrt{S(p_i) \cdot S(p_j)}}$ or $\frac{I}{\sqrt{H(p_i) \cdot H(p_j)}}$. With these choices, division of labor will be normalized to account for increases in numbers of individuals. The one thing that the matrix version of Simpson's index has against it is that it lacks the firm theoretical underpinnings that exist with the matrix version of Shannon's index. Otherwise, we envision the two normalized matrix indices as being equally useful measures of division of labor.

## Discussion

The problem with most existing division of labor statistics is that they examine only one individual or one task at a time and therefore cannot simultaneously account for the entire ensemble of interactions between individuals and tasks. This problem arises from building statistics whose inputs are probability vectors, with one vector per individual or task. This is not to say that these vector-based statistics may not be useful at capturing essential information about a given individual or task. These vector-based statistics are also based on well-accepted theory, relying on Shannon's index, Simpson's index, geometric mean, standard deviation, average deviation, and coefficient of variation – all of which have great statistical pedigrees.

The indices presented here are all distribution-free, hence there are no assumptions on the data other than the inputs being probabilities and, in some instances, that we not divide by zero. The assumption of no division by zero is not very restrictive: the only way this can happen is if each of the probabilities is identical to one another (i.e. equal amount of time spent on each task).

In most instances, a scalar output value makes for the best division of labor statistic. Sometimes it is useful examining a pair of division of labor statistics, such as simultaneously computing mutual entropy divided by marginal entropy of the rows and mutual entropy of the columns. In fact, we have used such ordered pairs of normalized mutual entropy to quantify whether breeding systems are monogamous, polygynous, or polyandrous (Gorelick et al., 2004). However, it is impossible to say whether one ordered pair or another has greater division of labor. Only scalar values can be rank ordered. The huge advantage of choosing a scalar measure of division of labor is that it allows for comparison across disparate data sets.

We acknowledge that no single number will ever capture all the facets of what is heuristically termed division of labor, such as task specialization, time spent performing each task, the probability of a task being performed, and task repertoire size. For example, when using normalized mutual entropy, the following six matrices all have zero division of labor: $\begin{pmatrix} 5 & 5 \\ 5 & 5 \end{pmatrix}$,

$\begin{pmatrix} 500 & 500 \\ 500 & 500 \end{pmatrix}$, $\begin{pmatrix} 10 & 0 \\ 10 & 0 \end{pmatrix}$, $\begin{pmatrix} 500 & 0 \\ 500 & 0 \end{pmatrix}$, $\begin{pmatrix} 5 & 5 & 5 & 5 \\ 5 & 5 & 5 & 5 \\ 5 & 5 & 5 & 5 \\ 5 & 5 & 5 & 5 \end{pmatrix}$,

and $\begin{pmatrix} 5 & 0 & 0 & 0 \\ 5 & 0 & 0 & 0 \\ 5 & 0 & 0 & 0 \\ 5 & 0 & 0 & 0 \end{pmatrix}$. The first two matrices each contains

a pair of generalists, who are either relatively lazy or very hard working. The third and fourth matrices each contain a pair of specialists, but these are degenerate cases insofar as there is only one task being performed by the entire population. The last two matrices show similar phenomena, but with populations of four individuals. Thus, one may want to incorporate the total number of individuals working, the total number of tasks performed, or the total amount of work performed into a division of labor statistic. Most of the statistics presented herein eschew these details, but in so doing, allow for comparisons of division of labor as the number of individuals, tasks, or total production varies. The fourteen statistics presented herein all have scalar outputs and capture a more subtle effect than mere increase in number of individuals, number of tasks or total production.

Having argued for a scalar output for a suitably general division of labor statistic, what do we want for the input to such a statistic? We argue that the key to a good division of labor statistic is to take the entire matrix of individuals and tasks as the input. With matrix inputs, there is less loss of information regarding interaction of rows and columns (individuals and tasks). Can the vector-based statistics we have enumerated be generalized so as to capture interactions between all individuals and tasks? No. There are higher-dimensional varieties of standard and average deviation. The standard or average deviation of a vector is a square symmetric matrix. This does not help in that – for a division of labor statistic – we are looking for a function whose range has smaller dimension than its domain, and not *vice versa*. There is therefore no apparent way to generalize division of labor indices based on standard or average deviation to simultaneously account for the ensemble of interactions that occurs amongst all individuals and tasks. Likewise, we know of no way to generalize the geometric mean to have a matrix input and scalar output. Thus, none of these approaches seem suitable for a division of labor statistic.

By contrast, Shannon's and Simpson's indices can be generalized to be matrix-based functions. Shannon himself did this in creating mutual entropy. The only short-falling of simply using mutual entropy as a measure of division of labor is that it is not normalized. Therefore, we recommend dividing mutual entropy by marginal entropy of either individuals or tasks (Colwell and Futuyma, 1987; Gorelick et al., 2004). We have also generalized Simpson's index to have a matrix input and scalar output. Again, the short-falling is that our generalized Simpson's index is not normalized, and so we divide the summand by whichever is larger, the number of individuals or number of tasks. These normalized matrix-input generalizations of Shannon's and Simpson's index (Indices 1c and 2f) should be the indices of choice when one wants to simultaneously examine division of labor amongst all individuals in a population.

## Acknowledgments

## References

Beshers S.N. and Traniello J.F.A. 1996. Polyethism and the adaptiveness of worker size variation in the attine ant *Trachymyrmex septentrionalius*. *J. Insect Behav.* **9**: 61 – 83

Boltzmann L. 1872. Weitere Studien über das Wärmegleichgewicht unter Gasmolekülen. *S. K. Akad. Wiss. Wien.* **66**: 275 – 370

Bormann, F.H. and G.E. Likens. 1967. Nutrient cycling. *Science* **155**: 424 – 429

Buckland S.T., Magurran A.E., Green R.E. and Fewster R.M. 2005. Monitoring change in biodiversity through composite indices. *Phil. Trans. R. Soc. London, Ser. B* **360**: 243 – 254

Colwell R.K. and Futuyma D.J. 1971. On the measurement of niche breadth and overlap. *Ecology* **52**: 567 – 576

Duncan O.D. and Duncan B. 1955. A methodological analysis of segregation indexes. *Am. Sociolog. Rev.* **20**: 210 – 217

Ernoult A., Bureau F. and Poudevigne I. 2003. Patterns of organisation in changing landscapes: implications for the management of biodiversity. *Landsc. Ecol.* **18**: 239 – 251

Gautrais J., Theraulaz G., Deneubourg J.-L. and Anderson C. 2002. Emergent polyethism as a consequence of increased colony size in insect societies. *J. Theor. Biol.* **215**: 363 – 373

Gibbs J.P. and Martin W.T. 1962. Urbanization, technology, and division of labor. *Am. Sociolog. Rev.* **27**: 667 – 677

Gibbs J.P. and Poston D.L. 1975. The division of labor: conceptualization and related measures. *Soc. Forces* **53**: 468 – 476

Gini C. 1912. *Variabilitá e mutabilitá*. Studi Economicoaguridici della facotta di Giurisprudenza dell Universite di Cagliari III, Parte II

Gorelick R. 2006. Combining richness and abundance into a single diversity index using matrix analogues of Shannon's and Simpson's indices. *Ecography* **29**: 525 – 530

Gorelick R., Bertram S.M., Killeen P. and Fewell J.H. 2004. Normalized mutual entropy in biology: quantifying division of labor. *Am. Nat.* **164**: 677 – 682

Julian G.E. and Cahan S. 1999. Undertaking specialization in the desert leaf-cutter ant *Acromyrmex versicolor*. *Anim. Behav.* **58**: 437 – 442

Keon T.L. and Carter N.M. 1985. Toward a clarification of the division of labor construct. *Hum. Relat.* **38**: 1131 – 1158

Keylock C.J. 2005. Simpson diversity and the Shannon-Wiener index as a special case of a generalized entropy. *Oikos* **109**: 203 – 207

Kolmes S.A. 1985. An information-theory analysis of task specialization among worker honey bees performing hive duties. *Anim. Behav.* **33**: 181 – 187

Labovitz S. and Gibbs J.P. 1964. Technology and division of labor. *Pac. Sociolog. Rev.* **7**: 3 – 9

Landau H.G. 1951. On dominance relations and the structure of animal societies. I. Effect of inherent characteristics. *Bull. Math. Biophys.* **13**: 1 – 19

Leontief W.W. 1951. Input-output economics. *Sci. Am.* **185**: 15 – 21

Martin W.T. and Gray L.N. 1971. Measurement of relative variance: sociological examples. *Am. Sociolog. Rev.* **36**: 496 – 502

Martínez I., Gil M.A. and López M.T. 1992. Analysis of mutual information measures in cluster sampling. *Appl. Math. Comp.* **52**: 389 – 402

Michener C.D. 1974. *The Social Behavior of the Bees*. Cambridge, Harvard University Press. 404 pp

Rao C.R. 1982. Diversity and dissimilarity coefficients: a unified approach. *Theor. Pop. Biol.* **21**: 24 – 43

Reeve H.K. and Keller L. 1995. Partitioning of reproduction in mother-daughter versus sibling associations: a test of optimal skew theory. *Am. Nat.* **145**: 119 – 132

Ricotta C. 2005. Through the jungle of biological diversity. *Acta Biotheor.* **53**: 29 – 38

Roughgarden, J. 2004. *Evolution's Rainbow*. Berkeley, University of California Press. 474 pp

Rushing W.A. 1968. Hardness of material as related to division of labor in manufacturing industries. *Admin. Sci. Q.* **13**: 229 – 245

Rushing W.A. and Davies V. 1970. Note on the mathematical formalization of a measure of division of labor. *Soc. Forces* **48**: 394 – 396

Seeley T.D. 1982. Adaptive significance of the age polyethism schedule in honeybee colonies. *Behav. Ecol. Sociobiol.* **11**: 287 – 293

Shannon C.E. 1948. A mathematical theory of communication. *Bell Sys. Tech. J.* **27**: 379 – 423, 623 – 656

Simpson E.H. 1949. Measurement of diversity. *Nature* **163**: 688

Smith A. 1776. *An Inquiry into the Nature and Causes of the Wealth of Nations*. Dublin, Whitestone. 590 pp

Smith D.L. and Snow R.E. 1976. The division of labor: conceptual and methodological issue. *Soc. Forces* **55**: 520 – 528

de Vries H. 1995. An improved test of linearity in dominance hierarchies containing unknown or tied relationships. *Anim. Behav.* **50**: 1375 – 1389

To access this journal online:
http://www.birkhauser.ch/IS