

# Multi-Level Selection and the Explanatory Value of Mathematical Decompositions

Christopher Clarke

*British Journal for the Philosophy of Science*\*

## Abstract

Do multi-level selection explanations of the evolution of social traits deepen the understanding provided by single-level explanations? Central to the former is a mathematical theorem, the multi-level Price decomposition. I build a framework through which to understand the explanatory role of such non-empirical decompositions in scientific practice. Applying this general framework to the present case places two tasks on the agenda. The first task is to distinguish the various ways of suppressing within-collective variation in fitness, and moreover to evaluate their biological interest. I distinguish four such ways: increasing retaliatory capacity, homogenising assortment, and collapsing either fitness structure or character distribution to a mean value. The second task is to discover whether the third term of the Price decomposition measures the effect of any of these hypothetical interventions. On this basis I argue that the multi-level Price decomposition has explanatory value primarily when the sharing-out of collective resources is 'subtractable'. Thus its value is more circumscribed than its champions Sober and Wilson ([1998]) suppose.

- 1 *Single-Level and Multi-Level Selection*
- 2 *Three Conditions on Explanatory Decompositions*
- 3 *The Multi-Level Price Decomposition*
- 4 *The Biological Interest Problem for Sober and Wilson*
- 5 *Explanatory Depth Whenever Resources are Subtractable*

---

\*Forthcoming 2014/2015. This is manuscript was created May 12, 2014— The content is identical to the copy submitted for final proofing.

- 6 *Other Cases, Alterations, and Roles?*
- 7 *The Second Term Doesn't Measure Between-Collective Variation*
- 8 *Alternative Approaches to Explanatory Depth*
- 9 *Conclusion*

## 1 Single-Level and Multi-Level Selection

One of the key variables in evolutionary theory is character–fitness covariance; the degree to which those organisms that possess a given character are statistically more likely to be fitter than those organisms that don't possess the character. Take for example a lion's inclination to hunt socially rather than on its own. Suppose that the fitness of each lion in a population is given by Table 1. So it's determined by whether or not that lion has this inclination to hunt socially, and by whether or not the lions that it interacts with have this inclination. Making

	who interacts with social hunters	who interacts with lone hunters
Fitness of social hunter	4	0
Fitness of lone hunter	1	1

Table 1: Example Fitness Matrix

some simple assumptions one can calculate that the covariance between character and fitness in this case is  $f_0(1 - f_0)(4f_0 - 1)$ ; where  $f_0$  is the proportion of the lion population who are presently social hunters.<sup>1</sup> Consider the case in which the population is evenly divided at present between social hunters and lone hunters; in other words  $f_0 = \frac{1}{2}$ . In these circumstances it follows that there is a positive covariance between social hunting and fitness, namely of  $\frac{1}{4}$ . This fact about covariance is key because it can provide a simple explanation of why the frequency of social hunters increased from the present generation of lions to the next generation. Explanation: lions inclined to hunt socially were—in the circumstances above—more likely to be fitter and this caused such lions to have relatively more offspring, most of whom inherited this inclination. And so the frequency of social hunters increased.

For reasons that will soon become clear I will call such explanations 'single-level selection' explanations. Such explanations are underwritten by the Robertson–Price identity. This equation describes how the covariance of character and

---

<sup>1</sup>Assume that lions form pairs of lions completely at random.

fitness determines the increased prevalence of a character in a population (Robertson [1966]; Price [1970]). This equation follows deductively from some common simplifying assumptions: that there is no migration into or out of the population; that the character in question is heritable and inherited without ‘transmission bias’; and that there are no stochastic effects at work (Price [1972]; Sober [1984]; Okasha [2006]). In the wake of Darwin’s *On the Origin of Species* single-level selection explanations have become so commonplace in evolutionary biology as to be unremarkable:

It would be advantageous to the *Melipona* [bee], if she were to make her cells closer together, and more regular in every way than at present; for then, as we have seen, the spherical surfaces would wholly disappear, and would all be replaced by plane surfaces; and the *Melipona* would make a comb as perfect as that of the hive-bee. . . . Thus, as I believe, the most wonderful of all known instincts, that of the hive-bee, can be explained by natural selection . . . (Darwin [1859], pp. 174–75)

Moving from explanations of concrete biological cases over to abstract mathematical models, this ‘single-level’ emphasis upon character–fitness covariance remains commonplace.<sup>2</sup> For example in textbook treatments of evolutionary theory one sees fitness matrices (such as Table 1) being used to identify the circumstances under which this character–fitness covariance will be positive, negative or zero (McElreath and Boyd [2007], p. 203). In the lion hunting case, for example, this depends upon the initial frequency  $f_0$  of social hunters in the population. Indeed one could describe the search in evolutionary game theory for so-called evolutionary stable states or strategies roughly as the search for the conditions under which character–fitness covariance is zero:  $f_0 = 0$ , or  $\frac{1}{4}$ , or 1 in this example.<sup>3</sup>

This illustrates how the covariance of character with fitness across the whole population is a central explanatory variable. Now in its multi-level form<sup>4</sup> the so-called Price equation decomposes this central variable into the sum of two other variables (Okasha [2006]). To put it briefly, one of these variables is supposed to relate in some sense to ‘selection at the level of individual lions’ and the other to ‘selection at the level of groups of lions’. I will say much more about these two variables in Section 3. For now it will suffice to say that both these variables

---

<sup>2</sup>McElreath and Boyd ([2007], §5.1) call the use of single-level explanations the ‘personal fitness approach’ to evolution.

<sup>3</sup>Note that this is a necessary but not sufficient condition for a distribution of characters across a population to constitute an evolutionarily stable distribution.

<sup>4</sup>The multi-level Price equation is a variation of the Price equation (Price [1972]), which itself is a more general form of the Robertson–Price identity (Robertson [1966]; Price [1970]).

are statistical functions of the distribution of character and fitness among lion groups.

Consider for example those cases in which selection for social hunting at the level of lion groups outweighed selection against social hunting at the level of individual lions. (Much more on this in Section 3.) In such cases the multi-level Price decomposition suggests a controversial explanation for the increase in the prevalence of social hunters from one generation to the next: group-level selection for social hunting outweighed individual-level selection against social hunting. As a consequence explanations that employ these two variables from the multi-level Price decomposition are often called ‘multi-level selection’ explanations.

The main focus of this paper will be the contrast between multi-level selection explanations and single-level selection explanations. And this will leave no time to say anything about the explanations afforded by selfish-gene theory (Dawkins [1976], [1982]) or inclusive fitness theory (Hamilton [1964]; Frank [1998]). Moreover considerations of space prevent me from discussing the alternative form of multi-level selection theory based on contextual analysis (Heisler and Damuth [1987]; Goodnight et al. [1992]) rather than the multi-level Price decomposition.

In contrasting the multi-level explanatory framework with the single-level framework I do not mean to imply that these frameworks offer competing explanations. As I define the concept, two explanations of the same case compete exactly when it is highly implausible, if not impossible, that they both be correct: for instance the explanation that the CIA shot Kennedy, and the explanation that Soviet agents shot Kennedy. In fact I’m happy to accept the so-called pluralist idea that multi-level explanations and single-level explanations—and for that matter selfish-gene and inclusive-fitness explanations—often posit the same process (Kerr and Godfrey-Smith [2002]) and so each framework can plausibly provide a correct explanation of the same case.

Instead, by contrasting multi-level explanations with single-level explanations, what I aim to do is to address the issue of explanatory depth. For example an explanation of why a car accelerated that specifies the car’s mechanics or the psychology of its driver provides a deeper explanation than merely citing the fact that the accelerator pedal was pressed. This shows how an explanation can be deeper than another without competing with it. On the one hand Sober and Wilson ([1998]) think that explanations of the evolution of social characters that employ the multi-level Price decomposition are deeper than single-level explanations. But on the other hand there are those who disagree. Maynard-Smith disagrees because he finds multi-level explanations altogether dubious;<sup>5</sup> whereas

---

<sup>5</sup>See Okasha ([2005], pp. 1000, 1004) for references and discussion of the complexities of

Dugatkin and Reeve ([1994], pp. 121, 124) disagree because they think multi-level explanations are fully equivalent to single-level explanations.<sup>6</sup>

The distinctive strategy of this paper will be to separate this issue of explanatory depth from the other issues in the ‘levels of selection’ literature with which it is entangled. In addressing it I will draw instead upon the general philosophical literature on explanation. Thus I will not discuss what it means for selection to ‘act at a particular level’ such as that of the group (Lloyd [1986], [2000]; Okasha [2006]), nor what it takes for something such as a group of organisms to count as a ‘biological individual’ (Clarke [Forthcoming]), nor whether groups can be vehicles in Dawkin’s ([1982]) sense, or interactors in Hull’s ([1981]) sense. Indeed one could perhaps think that there is no fact of the matter about such questions,<sup>7</sup> questions concerning vehicles or interactors say, but still think that there is a fact of the matter about the topic of this paper, namely the depth of the multi-level selection framework.

This focus on the explanatory depth of the multi-level Price decomposition will also raise wider philosophical questions. For the decomposition is a mathematical theorem: its truth isn’t contingent on what the world happens to be like; and one doesn’t need any scientific evidence to know that it’s true. Consequently one might wonder how such ‘non-empirical’ propositions could play a genuine role in scientific explanation (Pincock [2007]; Baker [2009]; Batterman [2010]). As Lange and Rosenberg ([2011], p. 593) point out in response to Sober ([2011]), it is ‘difficult to see how [propositions in evolutionary theory that are knowable a priori] could figure in causal explanations’.

So I will look beyond the philosophy of biology literature to explore how non-empirical decompositions such as the multi-level Price theorem can play an explanatory role. The suggestion will be—to put it somewhat laconically—that such decompositions highlight those constitutive relationships that help glue different factors in our explanatory reasoning together. Applying this suggestion to the multi-level Price decomposition shows that this decomposition has explanatory value, I will argue, primarily in cases in which the sharing-out of resources is ‘subtractable’. Thus the range of cases across which the decomposition provides deep explanations is more modest than its champions suppose.

---

Maynard Smith’s views.

<sup>6</sup>Things are not quite as clear cut as this. See Dugatkin and Reeve ([1994], p. 123). What is clear is that much confusion has been generated in contrasting multi-level selection explanations with their ‘individualist’ rivals, but not making clear what rivals one has in mind.

<sup>7</sup>See Sterelny ([1996]), Okasha ([2004a]), and Sarkar ([2008]) for discussion of this sort of pluralism.

## 2 Three Conditions on Explanatory Decompositions

What does one need to know in order to explain a phenomenon? In the philosophical literature a very popular suggestion is that one needs to know what would happen under certain ‘hypothetical alterations’ to the system in question. Would the phenomenon still have occurred if certain things had gone differently (Lewis [1986]; Woodward [2003])?<sup>8</sup> To explain why the economy shrank in 2008 for example it helps to know that the size of the economy would have been greater if banks had been more tightly regulated. So I am going to follow Lewis and Woodward in assuming that to explain is to answer important what-if-things-had-been-different questions. Accordingly the depth of an explanation is in proportion, roughly speaking, to the number of what-if questions it allows one to answer concerning important hypothetical alterations to the system in question. This measure of explanatory depth is by no means uncontroversial, but I will wait until Section 8 to examine it in further detail.

For the moment let’s note that this measure does not imply that answering each and every what-if question has explanatory value. After all, to explain why the economy shrank it does not help to know that the size of the economy would have been greater if extra-terrestrials had landed from outer space and donated a billion barrels of oil to the treasury. The what-if question about bank regulation is therefore different to the question about extra-terrestrial oil donation in that answering the former has explanatory value, but answering the latter does not. I will assume that the standard account of such differences is correct: we just happen to be more interested in hypothetical alterations to bank regulation than in far-fetched questions about extra-terrestrial oil donations.<sup>9</sup> The importance of a hypothetical alteration—and thereby its explanatory value—depends in this respect upon our personal interests.<sup>10</sup> (Accordingly, the notion of what is interesting to biologists will play a central role later in this paper.)

I will now use this Lewis–Woodward approach to explanation in order to build a toy model of how a non-empirical decomposition can play a modest role in explanation. Consider the following decomposition: the number of guests booked into a hotel is equal to the number of guests who are on holiday to ski plus the number of guests who are not on holiday to ski. This decomposition is non-empirical, guaranteed by the logical truth that everyone is either a skier or a non-skier. Compare this decomposition for example to a second decomposition,

---

<sup>8</sup>Lewis ([1986]) doesn’t put it in quite these terms. He says that to explain is to cite a cause; but for Lewis to cite a cause is just to say what could have gone differently such that the phenomenon wouldn’t have occurred.

<sup>9</sup>But see Hart and Honore ([1965]) for an alternative account.

<sup>10</sup>For an account that emphasises interests but not what-if questions see van Fraassen ([1977]) and Achinstein ([1983]).

a decomposition of the guests into those guests with blond hair and whose name begin with a 'K', on the one hand, and those guests who don't have both attributes on the other hand. The question I want to ask is this: when will a non-empirical decomposition (for example the first one) be more explanatorily valuable than any of the infinitely many other non-empirical decompositions that one might think of (such as the second decomposition)? To explore this question let us consider how the first decomposition fits into the following story.

(1) This winter has been unusually warm and so the average depth of snow on the Brixental ski slopes has been half a meter, in contrast to last winter's three meters. As a result (2a) there are two hundred skiers booked into the Brixental hotel, in contrast to last winter's nine hundred. (2b) And, like last year, there were one hundred non-skiers also booked into the Brixental hotel. Most of these non-skiers were there for the annual Wittgenstein conference. So applying our decomposition to (2a) and (2b) we see can see that (2) the hotel has had under five hundred guests rather than over five hundred as they did last winter. As a result the hotel has gone bankrupt.

Note that the low number of guests (factor 2) on its own provides a simple explanation of the bankruptcy. And this explanation is made deeper by adding the point about the lack of snow (factor 1). But to have a really satisfying explanation of the bankruptcy one needs also to be able to answer what-if questions of the following form: ( $Z$ ) what if  $x$  meters of snow had fallen, and other factors like the Wittgenstein conference been arranged in such-and-such a way?<sup>11</sup>

To answer such what-if questions one will typically reason as follows. 'In this hypothetical what-if scenario there would be  $g_a$  skiing guests on account of the snow; and there would be  $g_b$  non-skiing guests on account of the other factors such as the Wittgenstein conference. According to our decomposition this constitutes there being  $g$  guests in total. There would therefore be under five hundred guests, and so the hotel would be bankrupt. Alternatively: there would be over five hundred guests, and so the hotel would not be bankrupt.'

Let me break this down. To know our decomposition is to know a constitutive relationship  $X$ :  $g$  is constituted by  $g_a$  and  $g_b$ . And knowing this constitutive decomposition  $X$  is in practice how we come to know the following causal determination relationships  $Y$ : an interesting first factor (snowfall) combines with other factors (such as the Wittgenstein conference) to determine a second factor (total guests) which in turn determines the to-be-explained phenomenon (bankruptcy). And knowing these causal determination relationships  $Y$  in turn allows us to answer some important what-if questions  $Z$ . Thus this knowledge deepens our 'undecomposed' explanation (of the bankruptcy in terms of the total number of guests alone). In short our decomposition highlights a constitutive

---

<sup>11</sup> Arranged, for example, such that there were  $g_b$  non-skiing guests.

relationship that helps us to glue together the relevant factors in our explanatory reasoning.

In principle, of course, one could know these causal determination relationships  $Y$  without knowing the constitutive decomposition  $X$ . So the explanatory role of our decomposition is what one might call an ‘ancillary’ one. It is dispensable in principle, but not in practice.

It will be important for later to abstract three crucial elements from this toy example concerning the guests at the Brixental hotel.

*First element:* the value of a right-hand term in the decomposition ( $g_b$  non-skiing guests) is independent of the first factor (snowfall). After all, this term is a residual term that measures the effect of other factors only (such as the Wittgenstein conference). So its value is preserved by a hypothetical alteration to that first factor (eliminating snowfall). Observe that this element is crucial in that, without it, knowledge of the constitutive decomposition  $X$  will be of no help in calculating the causal dependencies  $Y$ . It will be very useful for later to note that this first element is equivalent to the following condition: the effect upon the left hand term ( $g$  total guests) of this alteration (eliminating snowfall) is measured by the attendant change to the value of a right-hand term in the decomposition ( $g_a$  skiing guests).

*Second element:* this hypothetical alteration (eliminating snowfall) is interesting. This element is crucial in that, without it, the what-if question  $Z$  would not be an important one. Hence answering this question would be of no explanatory value according to the Lewis–Woodward thesis about explanation; just as in my extraterrestrial oil donation example.

*Third element:* one knows how the value of the left-hand term ( $g$  total guests) determines the to-be-explained phenomenon (bankruptcy) in the circumstances. This element is crucial in that, without it, one could not use causal decomposition  $X$  to answer what-if question  $Z$ .

My conclusion is this: the Lewis–Woodward approach to explanation issues in three elements that are in general individually necessary and jointly sufficient for a non-empirical decomposition to provide explanatory value in the above manner. That is, to issue in an explanation of greater depth than an explanation (of the bankruptcy) in terms of only the left-hand term of a decomposition (the total number of guests).

I note in passing that the decomposition involving guests with blonde hair and names beginning with ‘K’ would in normal circumstances fail both criteria one and two.

I emphasise that the above are criteria only for the explanatory value of non-empirical decompositions, not empirical ones. To apply them to the case of empirical decompositions would be incorrect. For example the ideal gas law



$\ln(P) = \ln(V) + \ln(T)$  has clear explanatory value. But it fails my first criterion: when a gas is heated in an expandable chamber both the value of the  $\ln(V)$  term and of the  $\ln(T)$  term are altered as a result. So my first condition is not necessary as regards the explanatory depth of empirical decompositions, as opposed to non-empirical ones. Conversely the length of Edward Heath's premiership is equal to the length of Romano Prodi's premiership plus the length of John F. Kennedy's. This equation may well meet all my criteria, but it is too accidental to have any explanatory value.<sup>12</sup> So my three criteria are also not jointly sufficient as regards the explanatory value of empirical decompositions, as opposed to non-empirical ones.

At any rate the explanatory role played by the toy decomposition involving hotel guests, I will suggest, is the same explanatory role that many non-empirical decompositions play in the actual practice of science; in particular the multi-level Price decomposition in evolutionary biology.

### 3 The Multi-Level Price Decomposition

To spell out the multi-level Price decomposition let me introduce some standard formalism. Consider a population of individuals, be it a population of genes, cells, organisms or social groups; although the most intuitive case is when one takes individuals to be individual organisms. Take an arbitrary individual  $i$ . Let  $\omega_i$  denote that individual's (relative) fitness.<sup>13</sup> Let  $z_i$  denote the degree to which individual  $i$  possesses a particular character in which one is interested. This character will invariably be a 'social' character such as a lion's being inclined to hunt cooperatively or a vampire bat's being inclined to donate blood to other vampire bats who are in need. The multi-level Price decomposition states that:<sup>14</sup>

$$\text{Cov}(\omega, z) = \text{Cov}[\text{Exp}_g(\omega), \text{Exp}_g(z)] + \text{Exp}[\text{Cov}_g(\omega, z)] \quad (1)$$

What do these three terms mean? The left-hand term  $\text{Cov}(\omega, z)$  denotes the covariance of character with fitness across the whole population: to what extent do individuals who score high on character  $z$  tend statistically to be fitter than individuals in the population who score low on  $z$ ? For example are group hunters fitter on average than other lions?

Now imagine that our population of individuals is partitioned into collectives; so each individual is a member of exactly one collective. (I will leave it

<sup>12</sup>It certainly isn't invariant under interventions (Woodward [2003]). In contrast note that non-empirical decompositions are by definition maximally invariant.

<sup>13</sup>Relative fitness is defined to be an individual's absolute fitness divided by the mean fitness of all individuals in the population. I shall henceforth use 'fitness' to mean relative fitness.

<sup>14</sup>See Price ([1972]) and Hamilton ([1975]) for a seminal formulation. See Okasha ([2006]) for a very clear commentary.

entirely open what it is for an individual to be a member of a collective.) So  $\text{Cov}_g(\omega, z)$  denotes the covariance of character with fitness within collective  $g$ , rather than across the whole population: to what extent do individuals in the collective who score high on character  $z$  tend statistically to be fitter than those in the same collective who score low on character  $z$ ? Thus the third term of the decomposition  $\text{Exp}[\text{Cov}_g(\omega, z)]$  is an average of this measure across the whole population: on average do group hunters tend statistically to be fitter than those in the same collective who hunt alone?

Finally the second term.  $\text{Exp}_g(\omega)$  is the average fitness of the members of collective  $g$ . Let's call this the collective's fitness. Similarly  $\text{Exp}_g(z)$  is the average character of the members of collective  $g$ . Let's call this the collective's character.<sup>15</sup> So the second term of the multi-level Price decomposition  $\text{Cov}[\text{Exp}_g(\omega), \text{Exp}_g(z)]$  is the covariance between these two variables: to what extent do collectives who score high on character  $z$  tend statistically to be fitter than collectives who score low on character  $z$ ?<sup>16</sup> Putting this less technically and more intuitively: the second term of the decomposition measures the association *between collectives* of (collective) fitness with (collective) character, whereas the third term measures the association of (individual) fitness with (individual) character *within collectives*. Importantly the multi-level Price decomposition is a mathematical theorem, guaranteed by the logic of covariance and of expectation.

It is worth noting at this point that my third criterion for this decomposition to have explanatory value just requires that we know how the value of the left-hand term determines our to-be-explained phenomenon in the circumstances. And one does in this case. For one knows the Robertson–Price identity discussed in Section 1, which formally underwrites the intuition that the fitter character  $z$  is, so to speak, the more it will increase in frequency. So one knows how the value of the left-hand term (the degree of character–fitness covariance in the whole population) determines our to-be-explained phenomenon, the evolution of character  $z$ . So my third criterion is satisfied. Consequently, this paper will focus on the circumstances under which the multi-level Price decomposition satisfies my first and my second criterion for explanatory value.

#### 4 The Biological Interest Problem for Sober and Wilson

One suggested explanatory role for the multi-level Price decomposition emphasises the factor of within-collective variation (Sober and Wilson [1998]). And

---

<sup>15</sup>Thus I am focusing on what Damuth and Heisler ([1988]) call multi level selection type one, rather than type two.

<sup>16</sup>Strictly speaking the summation  $\text{Cov}[\ ]$  is over individuals in the population not collectives. So strictly speaking: to what extent do individuals that are part of collectives who score high on character  $z$  tend to be members of fit collectives?

by variation I strongly suspect that Sober and Wilson mean variation in fitness rather than variation in character.<sup>17</sup> Sober and Wilson’s key claim is that the third-term of the decomposition measures the effect of within-collective variation (pp. 32–33, 73–75). (Sober and Wilson also claim that the second-term of the multi-level Price decomposition measures the effect of between-collective variation. I will set the examination of this claim aside until Section 7.)

The general framework developed in Section 2 shows why Sober and Wilson’s key claim bears upon the explanatory depth of the multi-level Price decomposition. For this key claim is more or less an application of the first of my three criteria for explanatory value. Let  $\epsilon$  denote the effect of eliminating within-collective variation of fitness; in particular its effect upon character–fitness covariance across the whole population (as denoted by the left-hand term of the multi-level Price decomposition). First criterion: this effect  $\epsilon$  of this hypothetical alteration is measured by the attendant change in the value of a right-hand term in the decomposition (for example the third term). So Sober and Wilson’s key claim is more or less an application of the first of my three criteria for the multi-level Price decomposition to have explanatory value. Unfortunately Sober and Wilson do not provide an argument for this key claim. What follows is the most plausible way of developing such an argument in my view.

Take a population of individuals in an environment and consider the ‘fitness structure’ generated by that environment. This fitness structure is the mapping which specifies how an individual’s fitness is determined by her character and by the characters of the individuals with whom she interacts. Take for example the function  $\omega_i = 2\text{Exp}_g(z) - \frac{1}{2}z_i$ . Now consider a hypothetical alteration to this fitness structure such that each individual in any given collective  $g$  will now enjoy the same fitness as the other individuals in collective  $g$ . More precisely the fitness an individual is to enjoy under this alteration is identical to the mean fitness—prior to this alteration—of the individuals in her collective. For example in the above illustration  $\omega_i$  becomes equal to  $2\text{Exp}_g(z) - \frac{1}{2}\text{Exp}_g(z)$ . In other words it’s equal to  $\frac{3}{2}\text{Exp}_g(z)$ . Call this the ‘Structural Collapse to the Mean’ (SCM) alteration. This alteration is one straightforward way of eliminating any within-collective variation of individual fitness.

Note however that the SCM alteration preserves the mean fitness of the members of each collective, and thus preserves collective fitness. But individual character is also preserved; so collective character is preserved. Thus the SCM al-

---

<sup>17</sup>See pp. 54, 66–67, 80–91, 115, 139 of Sober and Wilson ([1998]) for textual evidence; indeed see Sober ([1984]). At any rate my criticism of Sober and Wilson’s idea as reconstructed in Section 6 will work just as well if you substitute ‘fitness’ for ‘character’ and ‘character’ for ‘fitness’. This is because covariance is symmetric:  $\text{Cov}(\omega, z) = \text{Cov}(z, \omega)$ . So the mathematical reasoning in my criticism will hold even if Sober and Wilson mean ‘variation in character’ rather than ‘variation in fitness’.

teration preserves the covariance of collective fitness with collective character. In other words it preserves the value of the second term of the multi-level Price decomposition. Let  $\epsilon$  denote the effect of SCM; in particular its effect upon character–fitness covariance across the whole population (as denoted by the left-hand term in the decomposition). SCM having preserved the value of the second term, it follows that this effect  $\epsilon$  is measured by the attendant change in the value of the third term in the decomposition. In other words my first criterion for explanatory value is satisfied here.<sup>18</sup>

Having established that my first criterion for explanatory value is satisfied with respect to hypothetical SCM alterations, can we now establish my second criterion? Is the SCM elimination of within-collective variation of fitness especially interesting to biologists? I will now argue that there is no general answer to this question: the answer depends very much on the biological details in each case.

Recall the example in which  $\omega_i = 2\text{Exp}_g(z) - \frac{1}{2}z_i$  which we can rewrite as  $2\text{Exp}_g(z) - z_i - \frac{1}{2}(-z_i)$ . Let's imagine that this describes the fitness structure for the *Polistes fuscatus* wasp in a given environment. Wasps with high  $z$  scores are hard workers. And wasps enjoy fitness benefits when they are in a collective whose members are hard working; hence the  $2\text{Exp}_g(z)$  term. But working hard requires a costly expenditure of energy; hence the  $-z_i$  term. But those lazy wasps who do not work hard run the risk of being stung by the queen, and indeed the risk of other forms of retaliation from the queen (Gamboa et al. [1990]); hence the  $-\frac{1}{2}(-z_i)$  term.

In the case of the *Polistes* wasp there is indeed a highly interesting way of altering the fitness structure that eliminates within-collective variation of fitness. One imagines an increase in retaliatory capacity: queens are better able to identify the lazy workers, or the queens increase the severity of the punishment for those who are so identified. In particular it will be interesting to know what would happen were the  $\frac{1}{2}$  coefficient—the retaliation parameter so to speak—to be altered such that each individual in a collective enjoys the same fitness; within-collective variation thus being eliminated. One can calculate that the answer is that the coefficient becomes 1 and that  $\omega_i$  becomes  $2\text{Exp}_g(z)$ .

It is crucial to note however that this highly interesting hypothetical alteration to fitness structure is distinct from the Structural Collapse to the Mean alteration I considered above. After all, recall that the SCM alteration has it instead that  $\omega_i$  becomes equal to  $\frac{3}{2}\text{Exp}_g(z)$ ; not to  $2\text{Exp}_g(z)$ . Furthermore there is nothing of especial biological interest I contend in the SCM alteration applied to our wasp population. For in this case a biologist has no reason to be interested

---

<sup>18</sup>Moreover one can easily show that SCM alters the value of the third term to zero. So the magnitude of this attendant change in the third term is given by the unaltered third term itself.

in hypothetical SCM alterations. Such alterations have no greater interest than hypothetical alterations that eliminate within-collective variation by letting  $\omega_i$  become  $\frac{7}{13}\text{Exp}_g(z)$ , or to  $\text{lnExp}_g(z)$ , or that collapse individual fitness to the collective median or the collective mode, or so on.

This illustrates how the Structural Collapse to the Mean alteration is not biologically interesting across every case in general. In other words SCM does not in general satisfy my second criterion for explanatory value. But I've been considering hypothetical SCM alterations in an attempt to develop Sober and Wilson's analysis into an argument that establishes a general explanatory role for the multi-level Price equation. And one can now see that this attempt has failed.

I emphasise that my intention here is not to criticise the application of the multi-level Price theorem to the *Polistes* wasp case. After all, the theorem is just a mathematical truth. Rather I am urging a more sanguine assessment of its explanatory value in this case. The decomposition doesn't obviously add any explanatory depth.

There will, of course, be some theorists who will resist my conclusion here by objecting to my relatively narrow conception of what is biologically interesting. I cannot hope to fully persuade such objectors. But I do hope to persuade them of a somewhat more modest point: the SCM alteration in the wasp case is just as interesting as the infinity of other hypothetical alterations to the distribution of fitnesses—such as those that let  $\omega_i$  become  $\frac{7}{13}\text{Exp}_g(z)$ , or  $\text{lnExp}_g(z)$ , or so on. It follows that, in the case of the *Polistes* wasp, the explanatory value of the multi-level Price decomposition will be just as great as the infinity of other mathematical decompositions of character–fitness covariance. So, the multi-level Price decomposition has no *special* explanatory value in the case of the *Polistes* wasp.

## 5 Explanatory Depth Whenever Resources are Subtractable

Two questions arise naturally from the last section. Firstly, still focusing on SCM alterations, can one develop Sober and Wilson's analysis into an argument that establishes the explanatory value of the multi-level Price equation in a more limited class of cases, rather than across all cases in general? And, secondly, can one appeal to hypothetical alterations other than SCM in order to establish an explanatory role for the decomposition, either generally or in a more limited range of cases?

Setting aside this latter question until Section 6, this section will tackle the former. It will identify a class of cases in which Structural Collapse to the Mean alterations are biologically interesting. In other words I identify a class of cases that satisfy my second criterion for explanatory value. These cases are, namely, those cases in which the sharing-out of resources amongst the individuals in a

collective is, in the parlance of economics, subtractable. But I've already shown in Section 3 that my third criterion is satisfied by the multi-level Price decomposition. And I've just shown in Section 4 that my first criterion is satisfied with respect to hypothetical SCM alterations. So all my three conditions are satisfied here. Thus this section establishes the explanatory value for the multi-level Price decomposition in a limited class of cases, namely those in which the sharing-out of resources is subtractable. Before getting down to business, I will need to invest a substantial amount of time illustrating what I mean by subtractability.

An excellent illustration of the subtractability of resources in a biological context is found in the literature on social or cooperative foraging (Giraldeau and Caraco [2000]). Many social foraging models can be thought of as having two parts. The resource acquisition part of the model describes how the amount of food that a collective of foragers will gather depends upon the cooperative behaviour of the members of the collective, and upon the environment. The resource sharing-out part of the model describes how the food that is foraged is divided amongst the individual members of the collective. Indeed there is an 'analytic separation' of the allocation of resources into a mechanism whereby a collective acquires its resources, and a mechanism whereby these resources are shared out amongst the individual members of the collective. I don't intend my point here to turn upon any substantial notion of 'mechanism'. Similarly I allow that two analytically separable mechanisms operate simultaneously, that they interact, and that they have overlapping parts. Instead, what I mean by 'analytic separation' is that there is a biologically interesting alteration of the manner in which resources are shared out amongst individuals, an alteration which leaves unaltered the manner in which resources are collectively acquired. To make this intuitive consider for example those 'scroungers' who have 'cheated' by refusing to cooperate during foraging. In many cases it is biologically interesting to ask what would occur if it became more difficult for scroungers to gain access to the food that the collective has foraged. What if, in the extreme, scroungers were excluded from these resources altogether? Thus I stipulate that properly-speaking resources are shared-out only if one can analytically separate resource allocation into a resource acquisition mechanism and a resource sharing-out mechanism. By this very token, resource sharing-out is subtractable only if these mechanisms are analytically separable. This is my first of two individually necessary and jointly sufficient conditions for subtractability.

Informally my second condition on subtractability is also rather intuitive: whenever one individual consumes a resource it must reduce the quantity of the resource available for other users to consume. To spell out the second condition formally I will make the simplifying assumption that one can use a single variable  $R_g$  to quantify the resources that a collective  $g$  happens to have acquired.

In a simple foraging case this is just the quantity of food that the collective has foraged. Furthermore I will assume that  $R_g$  is entirely determined by the ‘social’ character of each member of collective  $g$ , characters which one might represent by the vector  $\mathbf{z}_g$ . (In a simple foraging case this social character might measure how much energy the individual in question chooses to invest in the group hunt.) To emphasise this point I will often write collective resources  $R_g$  as  $R_g(\mathbf{z}_g)$  highlighting that it is a function of  $\mathbf{z}_g$ , and indeed of  $\mathbf{z}_g$  alone. Now consider the sum total of the fitnesses of the members of a collective  $g$ ; in formal terms  $\sum_G \omega_i$ . The sharing-out of collective resources is subtractable I stipulate only if this total fitness is entirely determined by collective resources  $R_g(\mathbf{z}_g)$ ; more specifically just in case this total fitness is an increasing function of collective resources. Choose the right scale on which to measure resources and this becomes the requirement that the fitness structure is characterised by:

$$\sum_G \omega_i = R_g(\mathbf{z}_g) \quad (2)$$

Why is this requirement on fitness structure fittingly described as a subtractability requirement? Notice that were any individual to be fitter than they actually are—but collective resources to remain as they actually are—then Equation 2 requires that some other individual or individuals would be less fit than they actually are, and by an equal amount. In the foraging case, holding fixed the amount of food collectively foraged, one individual’s gain in food/fitness is precisely counterbalanced by another’s loss.

It is of crucial importance to emphasise that the present requirement—concerning what would happen were collective resources to remain as they actually are—obviously does not entail that collective resources must remain as they actually are. Therefore many subtractable fitness structures will entail that collective resources vary according to the distribution of individual characters within the collective. In the foraging case for example the amount of food foraged  $R_g(\mathbf{z}_g)$  can vary depending on how the individuals cooperate during the hunt. So I emphasise that subtractability of resources does not entail that individuals are playing a zero-sum game that precludes them from cooperating to increase collective resources. A similar point: subtractability does not entail that the fitness structure in play is additive. In other words it does not entail that the fitness structure be given by  $\omega_i = \lambda z_i + \mu \text{Exp}_g(z)$ .

In summary, I stipulate that the sharing-out of resources is subtractable just in case (i) one can analytically separate resource allocation into a mechanism of resource acquisition by the collective and the mechanism that shares out these resources, and (ii) Equation 2 characterises the fitness structure in play.

A second illustration of the subtractability of resources comes from simple diploid genetics models. A genotype causes the organism in which it is instanti-

ated to exemplify a corresponding phenotype, and this organism interacts with the environment and has a number of offspring. And these offspring by extension are counted as the offspring of the genotype itself. Call this process the acquisition of a genotype's reproductive resources. (I'm happy to be fairly liberal about what counts as a resource.) Consider next that during meiosis each of the two alleles in the genotype will be copied to a certain number of gametes and so will enjoy a particular chance of being represented in each of the aforementioned organism's offspring. Call this the sharing-out of the genotype's reproductive resources amongst its two alleles. Again one can analytically separate resource allocation into collective resource acquisition and the sharing-out of these resources between individuals. For it is biologically interesting to ask what would occur if meiosis were to unfold differently: what if segregation distortion (Lyttle [1991]) occurred and the A-allele in an AB genotype enjoyed more than its fifty-percent share of reproductive resources (Maynard Smith and Szathmary [1995], §10)? So my first condition for subtractability is satisfied here. Equally my second condition for subtractability is also satisfied here: holding the genotype's resources fixed, an increased chance of the A-allele of being represented amongst the organism's offspring would be precisely counterbalanced by a decreased chance for the B-allele.

Finally, an example in which resources are *not* shared-out subtractably is that of the *Polistes* wasp. In this case a worker's fitness is sensitive to whether he is stung by the Queen. In virtue of this, avoiding being stung by the Queen is a key resource. But it would be absurd to attempt to analytically separate the allocation of this sting-avoidance resource into a mechanism whereby the collective acquires sting-avoidance, and a mechanism in which sting-avoidance is then shared out amongst individuals. So this resource is, by my definition, not 'shared out'. A second example in which resources are not shared-out subtractably is that of beavers building a channel from their dam to the river bank. I concede that one can analytically separate resource acquisition and resource sharing out here. But one beaver's using this channel does not exclude other beavers from doing likewise. So this sharing-out is not subtractable.

Almost there. I want now to make Equation 2 easier to work with mathematically. Consider the following constraint on the fitness  $\omega_i$  of each individual  $i$  in collective  $g$ :

$$\omega_i = \left(\frac{1}{n} - \alpha[z_i - \text{Exp}_g(z)]\right)R_g(\mathbf{z}_g) \quad (3)$$

Let me unpack this equation.  $\text{Exp}_g(z)$  is just the average character of the members of collective  $g$ . So  $[z_i - \text{Exp}_g(z)]$  denotes the degree to which our individual  $i$  scores especially highly on social character  $z$ . In other words whenever an individual has a perfectly average character then this becomes zero and the overall expression reduces to  $\frac{1}{n}R_g(\mathbf{z}_g)$ . In other words, whenever this is so, this indi-



vidual's fitness is equal to collective resources  $R_g(\mathbf{z}_g)$  divided by the number of members of the collective  $n$ . So whenever an individual is perfectly average she receives her 'fair share' of collective resources.

Similarly note that whenever an individual scores especially highly for 'social' character  $z$  then the  $-\alpha[z_i - \text{Exp}_g(z)]$  term will be negative; assuming  $\alpha$  is positive. So she will enjoy a lesser proportion of the collective's resources and thus she will be less fit. Conversely whenever an individual has an especially 'anti-social' character then this expression will be positive and so she will enjoy a greater proportion of collective resources and thus will be more fit. So the  $\alpha$  parameter denotes the degree to which 'anti-social' individuals can command an 'unfair' share of the resources that the collective has acquired. Thus parameter  $\alpha$  measures an important feature of the fitness structure generated by the environment, one that pertains to the sharing-out of resources between individuals as opposed to collective resource acquisition itself. (Table 2 illustrates the fitness structure that Equation 3 requires in a simple case; namely in the case of two-membered collectives, and in which an individual either has character  $z$  fully or not at all. In formal terms  $z = 0$  or  $z = 1$ .)

	who interact with a $Z$ individual	who interact with a non- $Z$ individual
Fitness of $Z$ individuals	$\frac{1}{2}R$	$(\frac{1}{2} - \frac{1}{2}\alpha)R'$
Fitness of non- $Z$ individuals	$(\frac{1}{2} + \frac{1}{2}\alpha)R''$	$\frac{1}{2}R''$

Table 2: Fitness of each individual in the subtractability case

Take the expression in round brackets in Equation 3 and sum it over all individuals in the collective. Since this necessarily sums to one it is evident that Equation 3 entails Equation 2. But I don't believe that to assume subtractability in the specific form of Equation 3 rather than more generally in the form of Equation 2 amounts to a significant loss in generality.<sup>19</sup> So from now on I will work with Equation 3 as part of my definition of subtractability, rather than Equation 2.

Having illustrated what I mean by subtractability, one can now get down to business. I will now argue that the multi-level Price decomposition has the ancillary role of answering questions about how character  $z$  would evolve if the environment were such that would-be anti-social individuals cannot gain unfair access to subtractable resources. Suppose that the sharing-out of resources amongst individuals is subtractable. Hence it can be characterised by a parameter  $\alpha$  which measures the degree to which the fitness-structure in play allows

---

<sup>19</sup>Frank's ([1995]) model however satisfies Equation 2 but not Equation 3.

anti-social individuals to access more than their fair share of collective resources. So intuitively, and as Equation 3 confirms, altering  $\alpha$  to become zero will reduce within-collective variation of fitness to zero. In these circumstances all individuals will receive an equal share of fitness, namely  $R_g(\mathbf{z}_g)$  divided by  $n$ . (One example of this would be an alteration of the visual environment such that would-be cheaters can be spotted, and thereby prevented from stealing extra resources.) So this hypothetical alteration of  $\alpha$  is a Structural Collapse to the Mean alteration. But I've already shown in Section 4 that all SCM alterations satisfy my first criterion for explanatory value: the effect  $\epsilon$  of this SCM alteration to  $\alpha$  will be measured by the attendant change to the value of the third-term in the multi-level Price decomposition.<sup>20</sup>

The second criterion for explanatory value requires that this alteration to  $\alpha$  be of interest to biologists. Note, however, that the existence of genuine sharing-out—by my definition—entails that one can analytically separate resource allocation into the acquisition of resources by the collective and the sharing-out of these resources amongst individuals. This in turn entails—again by my definition—that there is an interesting alteration to the mechanism of sharing-out resources amongst individuals, an alteration that does not alter how these resources were acquired by the collective. Therefore my definition of subtractable sharing-out guarantees that there will be biologically interesting alterations to  $\alpha$ . Here are two such cases. Case one:  $\alpha$  measures the degree to which visual environment is such that cheating foragers can go undetected, and therefore can steal resources rather than being excluded from them. Case two: alleles are taken as individuals, and genotypes are taken as collectives,  $\alpha$  measures the degree of so-called segregation distortion, the extent to which the meiotic environment allow selfish genetic elements to have more than their fair share of representation in the offspring organisms. These are just two examples of a biologically interesting  $\alpha$  parameter. So my second condition for explanatory value is (non-trivially)<sup>21</sup> satisfied.

But I've already shown in Section 3 that the third condition for explanatory

---

<sup>20</sup>Moreover one can show that the relationship between the third-term of the Price equation and  $\alpha$  is a linear one. For observe that it follows from Equation 3 that

$$\text{Cov}_g(\omega, z) = \text{Cov}_g\left(\left[\frac{1}{n} - \alpha z + \alpha \text{Exp}_g(z)\right]R_g, z\right) = \alpha R_g(\mathbf{z}_g)\text{Var}_g(z) \quad (4)$$

But one can substitute this into  $\text{Exp}[\text{Cov}_g(\omega, z)]$ , the third term of the multi-level Price decomposition, to yield  $\text{Exp}[\alpha R_g(\mathbf{z}_g)\text{Var}_g(z)]$ . And this yields  $\alpha \text{Exp}[R_g(\mathbf{z}_g)\text{Var}_g(z)]$ . For, being a feature of the environment,  $\alpha$  doesn't vary from collective to collective. So the third term of the Price decomposition depends linearly upon  $\alpha$ .

<sup>21</sup>More precisely: I've shown that all cases of sharing-out will by definition satisfy my second criterion for explanatory value. What these examples show is that the class of subtractable sharing-out cases is non-empty.

value is in general satisfied by the multi-level Price decomposition. So all three of my criteria are satisfied. Thus this section has established an explanatory role for the multi-level Price decomposition in a limited class of cases; namely cases in which the sharing-out of resources is subtractable. In such cases the multi-level Price decomposition deepens single-level explanations of the evolution of character  $z$  based on population-level character–fitness covariance alone. To put it intuitively, it has the ancillary role of answering questions about what would happen if the environment were such that anti-social individuals can no longer gain unfair access to subtractable resources.

## 6 Other Cases, Alterations, and Roles?

Section 4 showed that Sober and Wilson’s analysis cannot establish everything that Sober and Wilson want to establish. For it cannot establish the explanatory value of the multi-level Price decomposition across all cases in general: recall the case of retaliation in wasps. More precisely, Sober and Wilson’s analysis cannot achieve this when it is interpreted in terms of Structural Collapse to the Mean alterations. Instead Section 5 showed how Sober and Wilson’s analysis—interpreted in SCM terms—could be fleshed out to establish the explanatory value of the decomposition in the special case in which the sharing-out of resources is subtractable. Unfortunately I can’t see any other cases in which SCM alterations have any biological interest. (See my discussion in Section 4.) So I’m tentatively inclined to draw the following conclusion: the SCM alteration is only biologically interesting in cases in which collective resources are more or less subtractable. It follows that appealing to SCM alterations can establish no more than the explanatory value of the multi-level Price decomposition in cases in which resources are subtractable.

However, this doesn’t preclude there being other ways in which the multi-level Price decomposition might have explanatory value. Indeed there are hypothetical alterations other than SCM which eliminate within-collective variation. This naturally raises the following question: can one appeal to any of these other alterations in order to establish a further explanatory role for the multi-level Price decomposition? Perhaps the decomposition does indeed have a general explanatory role, or at very least a role in some cases in which resources are not subtractably shared-out. As I will illustrate momentarily, however, I can’t find any other hypothetical alterations which obviously satisfy my first and second criteria for explanatory value simultaneously; even for a limited range of cases. Therefore I’m inclined to jump to a further tentative conclusion: the multi-level Price decomposition plays an explanatory role only—or at least primarily—in cases in which resources are more or less subtractable. This section will sup-

port my claim here by examining the three obvious alternatives to the Structural Collapse to the Mean alteration: the ‘Increased Retaliatory Capacity’ alteration, the ‘Homogenizing Assortment’ alteration, and the ‘Character Collapse to the Mean’ alteration as I will label them.

*Character Collapse to the Mean.* Consider a collective of vampire bats composed of a few very fit members and many very unfit ones. Imagine for example a five-member collective containing individuals with fitnesses  $\omega = 1, 1, 1, 2,$  and  $10$ . Imagine altering the character of every member in the collective, and in turn their fitnesses, such that they are all moderately fit. Imagine in particular that this yields fitnesses of  $\omega = 3, 3, 3, 3,$  and  $3$ . By altering character, fitnesses have been collapsed to the collective mean. So within-collective variation in fitness, as measured by the third term of the multi-level Price decomposition, has been eliminated. Note that this Character Collapse to the Mean (CCM) alteration differs from the Structural Collapse to the Mean alteration in that it does not alter fitness via altering fitness structure; instead it does so by altering the frequency of the character in the population.

To see an immediate problem for appealing to CCM alterations, calculate the values of the second term in the multi-level Price decomposition for the example given in Table 3: the term is originally 90 but falls to 84 under the CCM alteration. Let  $\epsilon$  stand for the effect of CCM, in particular its effect upon char-

Original $z$	Original $\omega$	CCM $z$	CCM $\omega$
3	1	24	2
24	2	24	2
81	3	24	2
-	-	-	-
81	3	192	4
192	4	192	4
375	5	192	4

Table 3: Character Collapse to the Mean for two three-membered collectives and with  $\omega_i = \frac{1}{3} \sqrt[3]{z_i}$

acter–fitness covariance across the population (as denoted by the left-hand term of the multi-level Price decomposition). CCM having altered the value of the second term, it follows that this effect  $\epsilon$  is not measured by the attendant change to the third term of the multi-level Price decomposition. In other words, with respect to the CCM alteration, my first condition on explanatory value is not in general satisfied. Therefore one cannot appeal to the CCM alteration to identify a general explanatory role for the multi-level Price decomposition.

But this raises the following question: might appeals to CCM establish the explanatory value of the multi-level Price decomposition in a more limited range of cases, rather than across all cases in general? Now the only obvious non-gerrymandered range of cases that I can think of here is cases in which collective character maps one-to-one onto collective fitness. For one can show that the hypothetical CCM alteration does satisfy my first condition for explanatory value in such cases. This is because the CCM alteration will preserve collective fitness. And so, given the one-to-one mapping, it will preserve collective character. And so it will in turn preserve the covariance of collective fitness and collective character. In other words CCM will not alter the second term of the multi-level Price decomposition in this case. It follows that the attendant change to the third term does indeed measure CCM's effect in this case. So my first condition for explanatory value is met.

What about the second condition? I certainly do not deny that cases of one-to-one mapping constitute an interesting range of cases. For example this range of cases includes as a subset an important range of cases, namely those in which individual fitness is 'additive'.<sup>22</sup> (Additive cases are those in which fitness is a linear function of individual character and collective character:  $\omega_i = \lambda z_i + \mu \text{Exp}_g(z)$ . Thus collective character maps one-to-one onto collective fitness:  $\text{Exp}_g(\omega) = (\lambda + \mu) \text{Exp}_g(z)$ .) But I note incidentally that cases of one-to-one mapping excludes any form of synergism. In other words, it precludes individuals coordinating their activities so that the benefit to the collective is greater than the sum of each individual's own efforts.

This is all peripheral to the point I want to press here. More centrally, I question the biological interest of the Character Collapse to the Mean alteration itself. After all, the problems I identified in Section 4 with respect to the Structural Collapse to the Mean alteration can all be extended to Character Collapse to the Mean. For there is no *obvious* range of cases for which hypothetical collapses to the erstwhile mean are more biologically interesting than collapses to any other value (Section 4). Thus it is not obvious that CCM alterations ever satisfy my second requirement for explanatory value; even in a more limited range of cases.

*Homogenizing Assortment.* One biologically interesting alteration is the alteration to the mechanism of 'assortment', the mechanism that determines which individuals in a population join themselves into collectives with which other individuals. For example one might imagine that the mechanism of assortment is altered such that individuals only interact with individuals of a similar character. In the extreme case then assortment will be fully homogenous: within-collective variation in character and therefore within-collective variation in fit-

---

<sup>22</sup>See Birch ([2014]) for a discussion of assumptions similar to this additivity assumption but in a slightly different context.

ness will be zero. Thus the Homogenizing Assortment (HA) alteration differs from the CCM alteration in that it does not alter the overall composition of characters in the population, merely how individuals are assorted into collectives. It is clear that this HA alteration is in general biologically interesting. So it satisfies my second requirement for explanatory value.

Unfortunately, with respect to the Homogenizing Assortment alteration, my first criterion for explanatory value is never satisfied. To see this note that HA preserves the overall composition of characters in the population; by definition HA only alters how individuals in the whole population are grouped into collectives. So HA does not affect character–fitness covariance across the whole population (as denoted by the left-hand term of the multi-level Price decomposition). So trivially this (zero) effect of HA is not measured by the (non-zero)<sup>23</sup> attendant change to the value of the third term. In other words, with respect to the Homogenizing Assortment alteration, my first criterion on explanatory value is never satisfied. Therefore one cannot appeal to the Homogenizing Assortment alteration to identify any explanatory role for the multi-level Price decomposition at all.

*Increasing Retaliatory Capacity.* Recall the *Polistes* wasp example in which fitness was given by  $2\text{Exp}_g(z) - z_i - \frac{1}{2}(-z_i)$ . This is a special case of the more general fitness structure  $\omega_i = f(\mathbf{z}_g) - p(-z_i)$ ; where  $p$  is the parameter that measures retaliatory capacity (Section 4). Consider the hypothetical alteration in which this parameter is increased by  $\Delta p$ : queen wasps can for example more easily punish lazy workers, or punish them more severely. One can easily show that this Increasing Retaliatory Capacity (IRC) alteration increases the value of the second term of the multi-level Price decomposition, namely by  $\text{Var}[\text{Exp}_g(z)]\Delta p$ . Ruling out the trivial case in which there is no variation between collectives in collective character, this expression will be non-zero. In other words IRC alters the value of the second term of the decomposition. Let  $\epsilon$  stand for the effect of the IRC alteration, in particular its effect upon character–fitness covariance across the whole population (as denoted by the left-hand term of the multi-level Price decomposition). IRC having altered the second term, it follows that this effect  $\epsilon$  is never measured by the attendant change to the value of the third term in the decomposition.<sup>24</sup> So the IRC alteration fails my first criteria for the explanatory

---

<sup>23</sup>Homogenizing Assortment will eliminate the variation within any collective. So it will eliminate the character–fitness covariance within any collective. So it ensure that the value of the third term of the multi-level Price decomposition  $\text{Exp}[\text{Cov}_g(\omega, z)]$  will become zero. Setting aside the trivial case in which within-collective variation was already zero, this demonstrates that the attendant change to the value of the third term is non-zero.

<sup>24</sup>A similar point can be made about the second term. For the attendant change to the third term is, one can show:  $\text{Exp}[\text{Var}_g(z)]\Delta p$ . And this is only zero when there is no within-collective variation in individual character.

value of the decomposition in all but trivial cases.

To take stock: other than the Structural Collapse to the Mean alteration there are only three obvious hypothetical alterations to which one might appeal in order to establish the explanatory value of the multi-level Price decomposition. These alterations are Increasing Retaliatory Capacity, Homogenizing Assortment, and Character Collapse to the Mean. But I've shown decisively that one cannot appeal to the IRC or HA alterations to identify any explanatory role for the multi-level Price decomposition at all. And I've shown decisively that one cannot appeal to the CCM alteration to identify a general explanatory role for the multi-level Price decomposition at all. Moreover it's not *obvious* that we can find an explanatory role for CCM in any non-gerrymandered range of cases; for example those in which collective fitness maps one-to-one onto collective character. So one cannot obviously appeal to any of these three alterations—CCM, IRC, or HA—to establish any explanatory role for the multi-level Price decomposition. But the only other obvious alteration to which one might appeal is the SCM alteration, which I've already argued establishes the explanatory value primarily in cases in which the sharing-out of resources is more or less subtractable. Therefore this section lends some support to my tentative conclusion: the multi-level Price decomposition plays an explanatory role primarily in cases in which the sharing-out of resources is more or less subtractable.

## 7 The Second Term Doesn't Measure Between-Collective Variation

I've been examining Sober and Wilson's key idea that the third-term of the multi-level Price decomposition measures the effects of within-collective variation. But Sober and Wilson, I've already noted, also place a lot of weight upon an idea that is symmetrical to this one: the *second* term of the multi-level Price decomposition measures the effects of *between*-collective variation. But thus far I've ignored this symmetrical idea. This is because I think it is much more difficult to construct a plausible argument that favours it. The following is my best attempt, but one which ultimately fails.

Take a five-member collective with individual fitnesses of  $\omega = 1, 3, 6, 6,$  and  $9$ ; and thus of average fitness  $5$ . Consider a hypothetical alteration that changes the character of each member such that their fitness is 'boosted' by one unit, resulting in a five-member collective with fitnesses of  $\omega = 2, 4, 7, 7, 10,$  and thus of average fitness  $6$ . Note that it's a mathematical fact that this alteration won't alter within-collective variation of fitness. Consider also a second five-member collective with individual fitnesses of  $\omega = 1, 6, 8, 10,$  and  $10$ ; and thus of average fitness  $7$ . But this time consider a 'boost' of minus one unit, so that this second

collective now also has an average fitness of 6. Thus all collectives are altered to have the same collective fitness, in this case 6, thus eliminating between-collective variation in collective fitness. Consequently this Uniform Boosting alteration reduces to zero any covariance of collective fitness with other factors. Therefore  $\text{Cov}[\text{Exp}_g(\omega), \text{Exp}_g(z)]$ , the second term of the multi-level Price decomposition, will be zero.

However, let  $\epsilon$  denote the effect of this Uniform Boosting alteration, in particular its effect upon character–fitness covariance across the whole population (as denoted by the left-hand term in the multi-level Price decomposition). Unfortunately the second term of the multi-level Price decomposition does not in general measure this effect  $\epsilon$ . To see this, calculate the values of the third term in the multi-level Price decomposition for the example given in Table 4: the term is originally 62 but falls to 56 under the Uniform Boosting alteration. Uniform

Original $z$	Original $\omega$	Boost $z$	Boost $\omega$
3	1	24	2
24	2	81	3
81	3	192	4
-	-	-	-
81	3	24	2
192	4	81	3
375	5	192	4

Table 4: Uniform Boosting for two three-membered collectives and with  $\omega_i = \frac{1}{3} \sqrt[3]{z_i}$

Boosting having altered the value of the third term, it follows that this effect  $\epsilon$  is not measured by the attendant change to the second term of the multi-level Price decomposition. So, with respect to Uniform Boosting, the multi-level Price decomposition doesn't in general satisfy my first criterion for explanatory value.

One response might be to insist that nevertheless the attendant change in the second term measures effect  $\epsilon$  in a limited but non-gerrymandered class of cases. The only non-gerrymandered class of cases, however, for which this is obviously true are those in which an individual's fitness is a linear function of that individual's own character alone; put in formal terms  $\omega_i = mz_i + c$ . For whenever the fitness of each member of a collective is uniformly boosted by  $k$ , then each member's character will have been uniformly boosted by  $\frac{k}{m}$ , given this linear relationship. But the logic of covariance has it that  $\text{Cov}_g(\omega + k, z + \frac{k}{m}) = \text{Cov}_g(\omega, z)$ . So Uniform Boosting preserves the value of the third term in this case. It follows that this effect  $\epsilon$  is measured by the attendant change to the second term of the multi-level Price decomposition.



Unfortunately this class of cases is a completely irrelevant class for present purposes. For there's an intuitive sense in which there is no selection at the level of the collective at all in such cases. After all, in such cases individual fitness is not influenced by the collective at all. I have no doubt that Sober and Wilson would agree with this point. This is because, applying their own 1998 definition of 'trait groups', there are no genuine collectives in this special case. And hence there is no genuine collective-level selection.

So the problem remains: consider the effect  $\epsilon$  of eliminating between-collective variation via Uniform Boosting, in particular its effects upon character–fitness covariance across the whole population (as denoted by the left-hand term of the multi-level Price decomposition). Pace Sober and Wilson, there is no obvious class of cases for which this effect  $\epsilon$  is measured by the attendant change to the second term. So, with respect to the second term of the multi-level Price decomposition, there is no obvious class of cases for which my first criterion on explanatory value holds.

## 8 Alternative Approaches to Explanatory Depth

This paper has taken for granted that the depth of an explanation is in proportion, roughly speaking, to the number of important what-if questions that it allows one to answer. But why should one accept this? I cannot offer a full defence of this view, although interested evolutionary biologists might consult Woodward ([2003]), which has quickly become a philosophical classic. Instead this section will briefly examine the prospects for an alternative approach to explanatory depth, one that draws upon alternative accounts of explanation.

The first thing to note is that the philosophical literature contains scarcely any alternatives to the counterfactual account of explanatory depth. Why, for example, did the patient die? Hempel's Deductive Nomological approach might say that the following was a correct explanation: the patient ingested a large dose of digitalis, and it's a law that all people who ingest that dose will die soon afterwards (Hempel and Oppenheim [1948]). But Hempel's account is not an account of explanatory depth. For it does not offer us a criterion according to which this explanation counts as less deep than an explanation that includes details about how digitalis is metabolised and how it affects the heart. Hempel's approach is an account of explanatory correctness, not an account of the depth of a correct explanation.

Next consider Kitcher's ([1981]; [1989]) unificationist approach to explanation. Kitcher provides a criterion for what one might call explanatory promise, the ability of a candidate explanation to deepen one's understanding of what one already knows. And famously Kitcher's approach is a 'winner takes all' account.

Indeed it cannot be modified to admit degrees of explanatory promise on pain of admitting some embarrassing counter-examples (Woodward [2003], p. 368).<sup>25</sup> So—even if one were willing to equate explanatory promise with explanatory depth—Kitcher’s approach doesn’t delineate degrees of explanatory depth.

Kitcher’s approach should not be confused with the more modest—and thereby more plausible—idea that there are at least two virtues with respect to which an explanatory framework such as the multi-level selection framework can be assessed. The first virtue is what I’ve called depth, which I’ve urged is to be cashed out in terms of what-if questions. The second virtue is cashed out in terms of the framework’s scope of correct application: the broader the range of cases that can be correctly explained within that framework, the more ‘unifying’ the framework.<sup>26</sup> But it is evident that anyone tempted by this more modest unificationist idea will have no complaints with the assumptions that this paper has made about explanatory depth. All that the modest unificationist insists upon is that one also acknowledge the existence of an additional dimension to explanatory frameworks, unification qua broad scope of correct application.

Admittedly I’ve said very little about the relative scope of application of the single-level selection and multi-level selection frameworks. But this is because the answer is trivial: the multi-level selection framework has a narrower scope. After all, it embodies an extra restriction, namely that one’s population be partitioned into collectives. So, for this trivial reason, focusing on unification does not provide a sense in which multi-level selection explanations add value over and above single-level selection explanations.

Finally let’s consider the causal approach to explanation. Why have I been talking about the explanatory depth of the multi-level Price decomposition, rather than as Okasha ([2004b]; [2004c]) does of whether the decomposition is ‘causally adequate’ or ‘causally inadequate’? My main reason is that the notion of a decomposition’s being causally adequate is incredibly tricky (Okasha [Forthcoming]). That is why I have left the discussion in this paper incomplete as far as causal questions are concerned. But one might worry that in ‘ignoring’ causation the discussion in this paper is in danger of being not just incomplete but also unsound. I will now address this worry.

I’ve taken for granted through-out this paper that the depth of an explanation is, roughly speaking, in proportion to the number of important what-if questions that it helps to answer. And I’ve noted that the importance of a what-if question is in part determined by our personal interests. But philosophers who favour the causal approach to explanation might wish to place an additional restriction on

---

<sup>25</sup>Indeed see Woodward ([2003], §8) for what I take to be decisive counter-examples to the view overall.

<sup>26</sup>Birch ([2014], §5) proposes this more modest approach, although he seems to suggest that there is a sensible way of aggregating these two virtues into one overall score.

what counts as an important what-if question. The causal restriction: a what-if question is only important if the correct answer to it cites a cause of the to-be-explained event. I have no doubt that Lewis ([1986]), Lipton ([1991]), Ruben ([1990]), and Woodward ([2003]) amongst others would endorse this restriction.<sup>27</sup>

Adding this restriction, however, makes no difference to the soundness of the arguments of this paper. Firstly my criticism of Sober and Wilson in Sections 4 and 6 relied on the fact that certain what-if questions are uninteresting. And so my criticism required only that interestingness be a necessary condition for a what-if question to be important. It did not require that interestingness constitute the only necessary condition on importance. Secondly my positive point in Section 5 relied on the importance of questions about what would happen were parameter  $\alpha$  to be different. What happens to my argument if we add the requirement that  $\alpha$  has to be a cause of the evolution of social character  $z$  in order for  $\alpha$  to be part of the explanation for it? Nothing. For there is no reason to think that  $\alpha$ —an interesting feature of the environment that determines how much command anti-social individuals have over resources—cannot be a cause of the evolution of character  $z$ . So endorsing a causal approach to explanation does not generate a reason to resist the conclusions of this paper.

This concludes my defence of the measure of the depth of an explanation as, roughly, the number of important what-if questions that it helps to answer.

## 9 Conclusion

Sections 2 and 8 built and defended a general framework through which to understand the explanatory role of non-empirical decompositions such as the multi-level Price decomposition. Such decompositions have the ancillary role of describing the constitutive relationships that help glue different factors in our explanatory reasoning together. And I provided three individually necessary and jointly sufficient criteria for a non-empirical decomposition to play this role.

This framework highlighted two key questions for assessing whether multi-level selection explanations deepen the understanding provided by single-level selection explanations. Firstly, what hypothetical alterations are biologically interesting? What for example are the biologically interesting hypothetical alterations that eliminate within-collective variation in fitness? Secondly, does the third term of the Price equation measure the effects of any of these alterations upon character–fitness covariance across the whole population?

---

<sup>27</sup>But note that given Lewis' and Woodward's views of the nature of causation this restriction is a trivial one: roughly speaking, all answers to (the right sort of) what if things had been different questions cite causes.

I've considered four hypothetical alterations that eliminate within-collective variations in fitness: Increasing Retaliatory Capacity (Sections 4 and 6), Structural Collapse to the Mean (Sections 4 and 5), Homogenizing Assortment (Section 6), and Character Collapse to the Mean (Section 6). Only some of these hypothetical alterations are biologically interesting: HA is in general interesting; and SCM is interesting whenever resources are subtractable. And only some of these hypothetical alterations have their effects measured by the third term of the multi-level Price decomposition: the SCM alteration in all cases, and the CCM alteration in cases of one-to-one mapping of character to fitness.

Therefore none of these alterations are in general and simultaneously (a) biologically interesting and (b) such that their effects are measured by the attendant change to the third term of the multi-level Price decomposition. However, in the limited case in which resources are subtractable, the Structural Collapse to the Mean alteration is both biologically interesting and measured by the attendant change to the third term. The upshot is that the multi-level Price decomposition has explanatory value primarily when collective resources are more or less subtractable. Its value is more circumscribed than its champions Sober and Wilson ([1998]) believe.

Let me put the main thrust of the paper in intuitive form. What would happen if environmental conditions made it more difficult for anti-social individuals to access an unfair proportion of the subtractable resources acquired by their collective? I have argued that the explanatory value of the multi-level Price decomposition is that it helps us to answer such questions; questions about what would happen were the 'policing' of subtractable resources strengthened. But, I have shown, it does not help answer questions about other cases, or concerning other policing mechanisms such as retaliatory punishment or homogenizing assortment. This raises the question of how the paradigm policing mechanisms identified in Buss ([1987]) and Michod ([1999]) fit into my scheme for classifying policing mechanisms; and crucially whether these mechanisms issue in more or less subtractable resources.

## Acknowledgements

I am grateful to Jonathan Birch, Tim Lewens, Samir Okasha, Kim Sterelny, and two anonymous referees for their generous and helpful comments on the manuscript. European Research Council Grant agreement (284123) under the European Union's Seventh Framework Programme (FP7/2007-2013).

*Christopher Clarke*  
*Department of History and Philosophy of Science*

*University of Cambridge  
Free School Lane, Cambridge, CB2 3RH  
cjc84@cam.ac.uk*

## References

- Achinstein, P. [1983]: *The Nature of Explanation*, New York: Oxford University Press.
- Baker, A. [2009]: ‘Mathematical Explanations in Science’, *British Journal for the Philosophy of Science* 60: pp. 611–63.
- Batterman, R. W. [2010]: ‘On the Explanatory Role of Mathematics in Empirical Science’, *British Journal for the Philosophy of Science* 61: pp. 1–25.
- Birch, J. [2014]: ‘Hamilton’s Rule and its Discontents’, *British Journal for the Philosophy of Science* 65: pp. 381–411.
- Buss, L. W. [1987]: *The Evolution of Individuality*, Princeton NJ: Princeton University Press.
- Clarke, E. [Forthcoming]: ‘The Multiple Realizability of Biological Individuals’, *Journal of Philosophy*.
- Damuth, J. and I. L. Heisler [1988]: ‘Alternative Formulations of Multi-Level Selection’, *Biology and Philosophy* 3: pp. 407–30.
- Darwin, C. [1859]: *On the Origin of Species by Means of Natural Selection*: John Murray. Citations refer to the Revised Edition (2008) from Oxford University Press.
- Dawkins, R. [1976]: *The Selfish Gene*, Oxford: Oxford University Press.
- Dawkins, R. [1982]: *The Extended Phenotype*, Oxford: Oxford University Press.
- Dugatkin, L. A. and H. K. Reeve [1994]: ‘Behavioural Ecology and Levels of Selection: Dissolving the Group Selection Controversy’, in P. Slater, J. Rosenblatt, C. Snodown, and M. Milinski (eds), *Advances in the Study of Behaviour*, Volume 23: Academic Press, pp. 102–134.
- Frank, S. A. [1995]: ‘Mutual Policing and Repression of Competition in the Evolution of Co-operative Groups’, *Nature* 377: pp. 520–2.
- Frank, S. A. [1998]: *Foundations of Social Evolution*, Princeton NJ: Princeton University Press.
- Gamboa, G. J., T. L. Wacker, J. A. Scope, T. J. Cornell, and J. Shellman-Reeve [1990]: ‘The Mechanism of Queen Regulation of Foraging by Workers in Paper Wasps (*Polistes fuscatus*, Hymenoptera: Vespidae)’, *Ethology* 85: pp. 335–43.
- Giraldeau, L.-A. and T. Caraco [2000]: *Social Foraging Theory*, Princeton NJ: Princeton University Press.

- Goodnight, C. J., J. M. Schwartz, and L. Stevens [1992]: 'Contextual Analysis of Models of Group Selection, Soft Selection, Hard Selection, and the Evolution of Altruism', *American Naturalist* 140: pp. 743–61.
- Hamilton, W. D. [1964]: 'The Genetical Evolution of Social Behaviour', *Journal of Theoretical Biology* 7: pp. 1–16.
- Hamilton, W. D. [1975]: 'Innate Social Aptitudes in Man: An Approach from Evolutionary Genetics', in *Biosocial Anthropology*, New York: Wiley, pp. 133–55.
- Hart, H. L. A. and A. Honore [1965]: *Causation in the Law*, Oxford: Clarendon–OUP. Citations refer to the Second Edition (1985).
- Heisler, I. L. and J. Damuth [1987]: 'A Method for Analyzing Selection in Hierarchically Structured Populations', *American Naturalist* 130: pp. 582–602.
- Hempel, C. G. and P. Oppenheim [1948]: 'Studies in the Logic of Explanation', *Philosophy of Science* 15: pp. 135–175.
- Hull, D. L. [1981]: 'Units of Evolution: A Metaphysical Essay', in U. J. Jensen and R. Harré (eds), *The Philosophy of Evolution*, Brighton: Harvester Press, pp. 23–44.
- Kerr, B. and P. Godfrey-Smith [2002]: 'Individualist and Multi-Level Perspectives on Selection in Structured Populations', *Biology and Philosophy* 17: pp. 477–517.
- Kitcher, P. [1981]: 'Explanatory Unification', *Philosophy of Science* 48: pp. 507–531.
- Kitcher, P. [1989]: 'Explanatory Unification and the Causal Structure of the World', in P. Kitcher and W. Salmon (eds), *Scientific Explanation*, Minneapolis MN: University of Minnesota Press, pp. 410–505.
- Lange, M. and A. Rosenberg [2011]: 'Can There Be A Priori Causal Models of Natural Selection?', *Australasian Journal of Philosophy* 89: pp. 591–599.
- Lewis, D. K. [1986]: 'Causal Explanation', in *Philosophical Papers*, Volume 2, Oxford: Oxford University Press.
- Lipton, P. [1991]: *Inference to the Best Explanation*, London: Routledge.
- Lloyd, E. A. [1986]: 'Evaluation of Evidence in Group Selection Debates', *PSA: Proceedings of the Biennial Meeting of the Philosophy of Science Association* 1986: pp. 483–493.
- Lloyd, E. A. [2000]: 'Groups on Groups: Some Dynamics and Possible Resolution of the Units of Selection Debates in Evolutionary Biology', *Biology and Philosophy* 15: pp. 389–401.
- Lyttle, T. W. [1991]: 'Segregation Distorters', *Annual Review of Genetics* 25: pp. 511–57.

- Maynard Smith, J. and E. Szathmáry [1995]: *The Major Transitions in Evolution*, Oxford: Oxford University Press.
- McElreath, R. and R. Boyd [2007]: *Mathematical Models of Social Evolution: A Guide for the Perplexed*, Chicago IL: University of Chicago Press.
- Michod, R. E. [1999]: *Darwinian Dynamics: Evolutionary Transitions in Fitness and Individuality*, Princeton NJ: Princeton University Press.
- Okasha, S. [2004a]: ‘The “Averaging Fallacy” and the Levels of Selection’, *Biology and Philosophy* 19: pp. 167–84.
- Okasha, S. [2004b]: ‘Multi-Level Selection and the Partitioning of Covariance: A Comparison of Three Approaches’, *Evolution* 58: pp. 486–94.
- Okasha, S. [2004c]: ‘Multi-Level Selection, Covariance and Contextual Analysis’, *British Journal for the Philosophy of Science* 55: pp. 481–504.
- Okasha, S. [2005]: ‘Maynard Smith on the Levels of Selection Question’, *Biology and Philosophy* 20: pp. 989–1010.
- Okasha, S. [2006]: *Evolution and the Levels of Selection*: Oxford University Press.
- Okasha, S. [Forthcoming]: ‘The Relationship between Kin Selection and Multi-Level Selection’, *British Journal for the Philosophy of Science*.
- Pincock, C. [2007]: ‘A Role for Mathematics in the Physical Sciences’, *Nous* 41: pp. 253–275.
- Price, G. R. [1970]: ‘Selection and Covariance’, *Nature* 227: pp. 520–521.
- Price, G. R. [1972]: ‘Extension of Covariance Selection Mathematics’, *Annals of Human Genetics* 35: pp. 485–90.
- Robertson, A. [1966]: ‘A Mathematical Model of the Culling Process in Dairy Cattle’, *Animal Production* 8: pp. 95–108.
- Ruben, D.-H. [1990]: *Explaining Explanation*, London: Routledge.
- Sarkar, S. [2008]: ‘A Note on Frequency Dependence and the Levels/Units of Selection’, *Biology and Philosophy* 23: pp. 217–28.
- Sober, E. [1984]: *The Nature of Selection*, Chicago IL: Chicago University Press.
- Sober, E. [2011]: ‘A Priori Causal Models of Natural Selection’, *Australasian Journal of Philosophy* 89: pp. 571–589.



- Sober, E. and D. S. Wilson [1998]: *Unto Others: The Evolution and Psychology of Unselfish Behavior*, Cambridge MA: Harvard University Press.
- Sterelny, K. [1996]: 'The Return of the Group', *Philosophy of Science* 63: pp. 562–584.
- van Fraassen, B. C. [1977]: 'The Pragmatics of Explanation', *American Philosophical Quarterly* 14: pp. 143–150.
- Woodward, J. [2003]: *Making Things Happen: A Theory of Causal Explanation*, Oxford: Oxford University Press.