



**HAL**  
open science

## Matthew Effects via Team Semantics

Sabine Frittella, Giuseppe M Greco, Michele Piazzai, Nachoem M Wijnberg,  
Fan M Yang

► **To cite this version:**

Sabine Frittella, Giuseppe M Greco, Michele Piazzai, Nachoem M Wijnberg, Fan M Yang. Matthew Effects via Team Semantics. 2017. hal-01509419

**HAL Id: hal-01509419**

**<https://hal.science/hal-01509419>**

Preprint submitted on 17 Apr 2017

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Matthew Effects via Team Semantics

SABINE FRITTELLA, INSA Centre Val de Loire, U. of Orléans, LIFO EA 4022, France

GIUSEPPE GRECO, Department of Values, Technology and Innovation, Delft U. of Technology, the Netherlands

MICHELE PIAZZAI, Department of Values, Technology and Innovation, Delft U. of Technology, the Netherlands

NACHOEM M. WIJNBERG, Amsterdam Business School, U. of Amsterdam, the Netherlands; Department of Pure and Applied Mathematics, U. of Johannesburg, South Africa

FAN YANG, Department of Values, Technology and Innovation, Delft U. of Technology, the Netherlands

---

Social scientists call *Matthew effect* the self-reinforcing mechanisms whereby initially small advantages accrued by individuals, e.g. in reputation, capital, or access to opportunities, beget further advantage and result in growing inequality. While there is extensive literature on the Matthew effect, the notion has not been explicitly defined. In this paper, we take a first step in this direction by providing a formalisation of Matthew effects within the framework of team semantics, and by introducing a logic for analysing the properties of the dependence relations involved in Matthew effects. We also show via an example how to use this logic to formalise a statistical analysis of a Matthew effect.

CCS Concepts: • **General and reference** → **General conference proceedings**; • **Information systems** → **Relational database model**; • **Theory of computation** → **Logic and databases**; *Data modeling*; • **Mathematics of computing** → *Regression analysis*; • **Applied computing** → *Decision analysis*;

Additional Key Words and Phrases: Matthew effect; Team Semantics; Dependence Logic.

## ACM Reference format:

Sabine Frittella, Giuseppe Greco, Michele Piazzai, Nachoem M. Wijnberg, and Fan Yang. 2017. Matthew Effects via Team Semantics. 1, 1, Article 1 (April 2017), 13 pages.

DOI: 10.1145/nnnnnnn.nnnnnnn

---

## 1 INTRODUCTION

We provide a logical formalisation of the social phenomenon termed *Matthew effect* using team semantics. First introduced by Merton [12], the term Matthew effect refers to the self-reinforcing process whereby reputationally rich academics tend to get richer over time. In Merton’s words, it corresponds to “the accruing of large increments of peer recognition to scientists of great repute for particular contributions in contrast to the minimising or withholding of such recognition for scientists who have not yet made their mark” [13]. Outside the realm of academia, the Matthew effect has been invoked to explain positive feedbacks in e.g. the evaluation of athletes [10] and the conferral of public subsidies to firms [1]. Because of its ubiquity in social life, it has been recognized as a powerful engine of social, economic, and cultural inequality [14].

The role of Matthew effects is particularly evident in markets, where the uncertainty buyers face about the quality of products facilitates the consolidation of status hierarchies among sellers. This is because buyers consider the status of sellers as a good proxy for the quality of their products, so that higher status leads to greater visibility and access to resources, which help sellers achieve even greater status. For this reason, an extensive literature on the Matthew effect exists in sociology, economics, and management science. However, the Matthew effect is not precisely defined in this literature. As a result, researchers are hardly able to compare and integrate theoretical models and empirical findings. This motivates our present attempt to formalise the Matthew effect in logic.

The logical framework we propose for our formalisation is the framework of team semantics. Introduced originally by Hodges [8, 9] and later advanced by Väänänen [15], team semantics is a novel and effective logical tool for analysing the notion of *dependency* that is fundamental to the social and natural sciences. These dependencies usually manifest themselves only in the presence of multiple observations. Team semantics thus evaluates formulas under sets

---

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

© 2017 Copyright held by the owner/author(s). XXXX-XXXX/2017/4-ART1 \$15.00

DOI: 10.1145/nnnnnnn.nnnnnnn

of assignments (called *teams*), instead of single assignments as in the standard Tarskian semantics. Teams can be easily conceived as sets of rows in tables or data sets. The flexible and multidisciplinary interpretations of teams engendered a rapid development of logics based on team semantics in recent years. Notable such logics are *dependence logic* [15], *independence logic* [5] and *inclusion logic* [3, 4], which focus on characterising *functional dependence*, *independence* and *inclusion relations* among variables, respectively. In this paper we focus on the dependence relations relevant to Matthew effects. To the best of our knowledge, these have not been formalised in the literature. We take a first step in this direction by introducing a new logic based on team semantics, called *logic of Matthew effects* (ML), and analysing basic properties of these dependence relations within the framework of ML.

*Structure of the paper.* In Section 2, we define the syntax and semantics of our logic ML. In Section 3, we formalise four distinct types of Matthew effects and study their properties. In Section 4, we use ML to formalise an existing study on the Matthew effect in science [2].

## 2 LOGIC OF MATTHEW EFFECTS

We first present the general framework of the logic ML, including team semantics, mathematical definitions, and notation. We provide a brief presentation of time series regression, a common tool to analyse trends in dynamic data. For more details on this topic, the reader can refer to [6]. In Section 2.2, we introduce the basic dependence relations we use as building blocks to define Matthew effects. In Section 2.3, we present the syntax and semantics of ML.

### 2.1 General framework

We introduce the general framework to talk about data sets and regressions over data sets. In Appendix A.1, we present a toy example of a data set (Table 1) and a Matthew effect to exemplify the team semantics. We refer to this example throughout this subsection. In this paper, we assume that the empirical context where the Matthew effect occurs is described by a first-order model  $M$  of some signature  $\mathcal{L}$  that will be specified in the sequel. The model  $M$  is assumed to have a three-sorted domain: *data sort*, *regression sort*, and *duration sort*.

*Three-sorted domain.* Elements of the *data sort* are all possible values in the data sets, e.g. real numbers or names. We include the value undefined to the possible values for variables of type data sort. We use  $w, x, y, z, \dots$  (possibly with subscripts) to stand for variables of this sort. Data sets can be presented as tables (e.g. Table 1), and the variables of data sort can be understood as attributes in such tables. We assume that the time attribute is present in every data set, and we reserve the letter  $t$  for the time variable. In this setting, an observation in a data set or a row of a table corresponds to an assignment  $s$  that assigns to each variable  $x$  of data sort a value  $a$  of data sort in the domain of the underlying model  $M$ . Formally, a data set  $D$  consisting of several observations is a set of assignments, which we also call *team*. In this paper, the terms *data set* and *team* are used interchangeably. As we explained in the introduction, the logic we introduce for Matthew effects adopts team semantics, where formulas are evaluated on teams.

Elements of the *regression sort* are names of regressions, and we use  $r, r', \dots$  to stand for variables of this sort. For a given model  $M$ , we call *granularity* the minimal time interval  $\delta$ , between two observations. For instance, in Table 1 the granularity is  $\delta = 1 \text{ year}$ . Different regressions can be performed on any data set using different variables. Consecutive observations for a single variable can be aggregated over time intervals that correspond to particular multiples of the granularity  $\delta$ . The intervals of time (i.e., the natural numbers that correspond to the multiples of the minimal time interval) form the domain of the *duration sort*. We use  $\ell, \ell', \dots$  to stand for variables of this sort.

*Regression functions.* Each regression generates one *regression function* for each dependent variable  $y$  under consideration. A regression function can be represented informally as:

$$y_{(t)} = \sum_{\{i_1, \dots, i_k\} \in \mathcal{X}} \beta_{i_1 \dots i_k} (x_{i_1})_{(t-\delta)} \cdots (x_{i_k})_{(t-\delta)} + \epsilon, \quad (1)$$

where  $\{x_1, \dots, x_n\}$  is a set of independent variables,  $z_{(t)}$  is the value of variable  $z$  at time  $t$ ,  $\mathcal{X}$  is a non-empty downward closed subset of  $\mathcal{P}(\{1, \dots, n\})$  (i.e.,  $B \subseteq A \in \mathcal{X}$  implies  $B \in \mathcal{X}$ ), the positive real number  $l \cdot \delta$  is the time interval of the regression,  $l$  is a natural number,  $\epsilon$  is an error term, and each coefficient  $\beta_{i_1 \dots i_k}$  is a real number. Most studies in the social sciences, including those on Matthew effects, use regressions of degree two or less, for instance:

$$y_{(t)} = \alpha + \beta_1(x_1)_{t-\delta} + \beta_2(x_2)_{t-\delta} + \beta_3(x_3)_{t-\delta} + \beta_4(x_1)_{t-\delta}(x_2)_{t-\delta} + \beta_5(x_2)_{t-\delta}(x_3)_{t-\delta} + \epsilon. \quad (2)$$

In a data set  $D$ , if  $s$  and  $s'$  are two observations that respectively contain the data for time  $t$  and  $t-\delta$  (i.e.,  $s(t) = s'(t)+\delta$ ), then the above regression equation (2) can be formally represented as:

$$s(y) = \alpha + \beta_1 s'(x_1) + \beta_2 s'(x_2) + \beta_3 s'(x_3) + \beta_4 s'(x_1)s'(x_2) + \beta_5 s'(x_2)s'(x_3) + \epsilon.$$

If  $x_1$  and  $x_2$  are the focal variables, we will sometimes abbreviate the expression above as:

$$s(y) = \beta_1 s'(x_1) + \beta_2 s'(x_2) + \beta_4 s'(x_1)s'(x_2) + q(s'(x_1), s'(x_2), s'(x_3)).$$

*Analysed data set.* We call a data set  $D$  associated with regression  $R_1, \dots, R_k$  with durations of time  $l_1, \dots, l_k$ , respectively, an *analysed data set*, denoted  $(D, R_1, l_1, \dots, R_k, l_k)$ . Formally, we view this analysed data set as a data set extended from  $D$  by adding  $2k$  columns with attributes  $r_1, \dots, r_k, \ell_1, \dots, \ell_k$ , where each column  $r_i$  has a constant value  $R_i$ , and each column  $\ell_i$  has a constant value  $l_i$ . For simplicity, we denote this analysed data set also by  $D$ , and write  $D(r_i)$  for the unique value of the attribute  $r_i$  in  $D$ , and  $D(\ell_i)$  for the unique value of the attribute  $\ell_i$ . Technically, we may view a name  $R$  of a regression analysis as a function mapping each dependent variable  $y$  under consideration to a polynomial  $R(y)$  of the form (1).

## 2.2 Basic dependence relations.

We define basic dependence relations and atomic formulas that we later use to define Matthew effects.

*Definition 2.1 (Basic dependence relations and their team semantics).* Let  $x_1, \dots, x_n, y$  be data sort variables,  $r$  be a regression sort variable, and  $\ell$  be a duration sort variable.

- The atomic formula  $x_1, \dots, x_n \overset{r}{\mathcal{P}}_{\ell} y$  characterises the notion of  $y$  being **positively  $\ell$ -dependent on  $x_1, \dots, x_n$  with respect to the regression  $r$** . We say that the defined relation is true in a data set  $D$  with an underlying model  $M$ , denoted  $M \models_D x_1, \dots, x_n \overset{r}{\mathcal{P}}_{\ell} y$ , iff for all  $s, s' \in D$ ,

$$s(t) = s'(t) + D(\ell) \cdot \delta^M \implies s(y) = \beta_1 s'(x_1) + \dots + \beta_n s'(x_n) + q(s'(\bar{x}), s'(\bar{w})).$$

where each  $\beta_i$  is significantly greater than 0 (see Appendix B.2), and  $D(r)(y)$  is the polynomial represented above.

- The atomic formula  $x_1, \dots, x_n \overset{r}{\mathcal{N}}_{\ell} y$  characterises the notion of  $y$  being **negatively  $\ell$ -dependent on  $x_1, \dots, x_n$  with respect to  $r$** , and we define  $M \models_D x_1, \dots, x_n \overset{r}{\mathcal{N}}_{\ell} y$  the same way as above except that we now require each  $\beta_i$  to be significantly smaller than 0.
- The atomic formula  $x_1 \otimes \dots \otimes x_n \overset{r}{\mathcal{M}}_{\ell} y$  characterises the notion of  $y$  being **positively moderated  $\ell$ -dependent on  $x_1, \dots, x_n$  with respect to  $r$** , and we define  $M \models_D x_1 \otimes \dots \otimes x_n \overset{r}{\mathcal{M}}_{\ell} y$  iff for all  $s, s' \in D$ ,

$$s(t) = s'(t) + D(\ell) \cdot \delta^M \implies s(y) = \sum_{k=1}^n \sum_{1 \leq i_1 < \dots < i_k \leq n} \beta_{i_1 \dots i_k} s'(x_{i_1}) \dots s'(x_{i_k}) + q(s'(\bar{x}), s'(\bar{w})),$$

where  $\beta_{i_1 \dots i_k}$  is significantly greater than 0, and  $D(r)(y)$  is the polynomial represented as above.

- The atomic formula  $x_1 \otimes \dots \otimes x_n \overset{r}{\mathcal{NM}}_{\ell} y$  characterises the notion of  $y$  being **negatively moderated  $\ell$ -dependent on  $x_1, \dots, x_n$  with respect to  $r$** , and we define  $M \models_D x_1 \otimes \dots \otimes x_n \overset{r}{\mathcal{NM}}_{\ell} y$  the same way as above except that we now require  $\beta_{i_1 \dots i_k}$  to be significantly smaller than 0.

*Definition 2.2 (Atomic formulas for positive and negative dependency).* Let  $\mathbb{X}$  denote the set of strings  $X$  of dependent variables of the form

$$x_{11} \otimes \dots \otimes x_{1n_1}, \dots, x_{k1} \otimes \dots \otimes x_{kn_k}.$$

For every string  $X \in \mathbb{X}$  such that  $X = x_{11} \otimes \dots \otimes x_{1n_1}, \dots, x_{k1} \otimes \dots \otimes x_{kn_k}$  and for every duration of time  $\ell \cdot \delta$ , we introduce the atomic formulas  $X \overset{r}{\mathcal{P}}_{\ell} y$  and  $X \overset{r}{\mathcal{N}}_{\ell} y$ , defined as follows:

- $M \models_D x_{11} \otimes \dots \otimes x_{1n_1}, \dots, x_{k1} \otimes \dots \otimes x_{kn_k} \overset{r}{\mathcal{P}}_{\ell} y$  iff for all  $s, s' \in D$ ,

$$s(t) = s'(t) + D(\ell) \cdot \delta^M \implies s(y) = \sum_{i=1}^k \sum_{m=1}^{n_i} \sum_{1 \leq j_1 < \dots < j_m \leq n_i} \beta_{i, j_1 \dots j_m} s'(x_{i j_1}) \dots s'(x_{i j_m}) + q(s'(\bar{x}), s'(\bar{w})), \quad (3)$$

where each  $\beta_{i, 1, \dots, n_i}$  is significantly greater than 0, and  $D(r)(y)$  is the polynomial represented as above.

- $M \models_D x_{11} \otimes \dots \otimes x_{1n_1}, \dots, x_{k1} \otimes \dots \otimes x_{kn_k} \overset{r}{\mathcal{NM}}_{\ell} y$  iff the same as the above holds except that each  $\beta_{i, 1, \dots, n_i}$  is now required to be significantly smaller than 0.

### 2.3 The logic for Matthew effects

Here, we define the signature, the syntax, and the semantics of ML.

*Definition 2.3 (The signature  $\mathcal{L}$ ).* The signature  $\mathcal{L}$  of ML contains the constant functions 0 and  $\delta$ , the unary functions  $\{\text{Effect}_{y,X} \mid y \in \text{Var}, X \in \mathbb{X}\}$ , the unary predicate Small, and the binary predicates  $\gg, \ll, \text{ and } \approx$ .

The function symbols  $\text{Effect}_{y,X}$  are used to refer to certain combinations of the coefficients in the polynomials of the regression functions. For every variables  $x, y$ , every string  $X$  of dependent variables,  $\text{Effect}_{y,X}(r)$  is a regression sort argument, whose interpretation in the intended models  $M$  under a data set  $D$  is defined as:

- iff  $X = x_{11} \otimes \dots \otimes x_{1n_1}, \dots, x_{k1} \otimes \dots \otimes x_{kn_k}$ , and  $\beta_1 x_{11} \dots x_{1n_1}, \dots, \beta_k x_{k1} \dots x_{kn_k}$  are terms in the polynomial  $D(r)(y)$ , then for every  $s \in D$ ,  $\text{Effect}_{y,X}^M(s(r)) = \beta_1 + \dots + \beta_k$ . Otherwise, set  $\text{Effect}_{y,X}^M(s(r)) = \text{undefined}$ .

We equivalently denote  $\text{Effect}_{y,X}(r)$  by  $\text{Effect}(r, y, X)$ . The use of these terms becomes clearer in Section 4. The constant symbols 0 and  $\delta$  are to be interpreted as the natural number 0 and the granularity of the data sets  $\delta$  respectively. The predicate symbols Small,  $\ll, \gg$  and  $\approx$  are to be interpreted as “small,” “significantly smaller than,” “significantly greater than,” and “equivalent to,” respectively.

*Definition 2.4 (Syntax).* Let  $\text{Var}_0, \text{Var}_1$  and  $\text{Var}_2$  be respectively countable sets of variables of data sort, regression sort, and duration sort. The syntax of ML is defined as follows:

Terms of data sort  $\alpha ::= x \mid 0 \mid \delta \mid \text{Effect}_{y,X}(r)$

Terms of duration sort  $\beta ::= \ell$

Terms of regression sort  $\gamma ::= r$

Formulas  $\phi ::= X \overset{r}{\smile}_{\ell} y \mid X \overset{r}{\smile}_{\ell} y \mid \text{Small}(\alpha) \mid \alpha \ll \alpha \mid \alpha \gg \alpha \mid \alpha \approx \alpha \mid \phi \wedge \phi \mid \exists^1 x \phi \mid \exists^1 \ell \phi \mid \exists^1 r \phi$

where  $x \in \text{Var}_0, r \in \text{Var}_1$  and  $\ell \in \text{Var}_2$ .

*Definition 2.5 (Free variables).* For compound formulas, the sets of free variables of each sort are defined as usual. For atomic formulas, the sets of free variables of each sort are defined as follows:

- The set  $\text{Fv}_0(\phi)$  of free variables of data sort is defined as
  - for  $X = x_{11} \otimes \dots \otimes x_{1n_1}, \dots, x_{k1} \otimes \dots \otimes x_{kn_k}$ ,
$$\text{Fv}_0(X \overset{r}{\smile}_{\ell} y) = \text{Fv}_0(X \overset{r}{\smile}_{\ell} y) = \{x_{i,j} \mid 1 \leq i \leq k \text{ and } 1 \leq j \leq n_i\} \cup \{y, t\}^1,$$
  - $\text{Fv}_0(\text{Small}(\alpha)) = \text{Fv}_0(\alpha)$ ,
  - $\text{Fv}_0(\alpha \ll \beta) = \text{Fv}_0(\alpha \gg \beta) = \text{Fv}_0(\alpha \approx \beta) = \text{Fv}_0(\alpha) \cup \text{Fv}_0(\beta)$ , and
  - the set  $\text{Fv}_0(\alpha)$  is defined inductively as  $\text{Fv}_0(x) = \{x\}$  and  $\text{Fv}_0(0) = \text{Fv}_0(\delta) = \text{Fv}_0(\text{Effect}_{y,X}(r)) = \emptyset$ .
- The set  $\text{Fv}_1(\phi)$  of free variables of regression sort is defined as
  - $\text{Fv}_1(X \overset{r}{\smile}_{\ell} y) = \text{Fv}_1(X \overset{r}{\smile}_{\ell} y) = \{r\}$ ,
  - $\text{Fv}_1(\text{Small}(\alpha)) = \text{Fv}_1(\alpha)$  and  $\text{Fv}_1(\alpha \ll \beta) = \text{Fv}_1(\alpha \gg \beta) = \text{Fv}_1(\alpha \approx \beta) = \text{Fv}_1(\alpha) \cup \text{Fv}_1(\beta)$ , and
  - the set  $\text{Fv}_1(\alpha)$  is defined inductively as  $\text{Fv}_1(x) = \text{Fv}_1(0) = \text{Fv}_1(\delta) = \emptyset$  and  $\text{Fv}_1(\text{Effect}_{y,X}(r)) = \{r\}$ .
- The set  $\text{Fv}_2(\phi)$  of free variables of duration sort is defined as
  - $\text{Fv}_2(X \overset{r}{\smile}_{\ell} y) = \text{Fv}_2(X \overset{r}{\smile}_{\ell} y) = \{\ell\}$  and
  - $\text{Fv}_2(\phi) = \emptyset$  for any other atomic formula.

Formulas with sets  $\text{Var}_0, \text{Var}_1, \text{Var}_2$  of free variables of data sort, regression sort, and duration sort are evaluated on a model  $M$  with respect to *teams* over  $V_0 \cup V_1 \cup V_2$ , i.e., *sets*  $D$  of assignments  $s : V_0 \cup V_1 \cup V_2 \rightarrow M$ .

*Definition 2.6 (Semantics).* We define inductively the satisfaction relation  $M \models_D \phi$  as follows:

- See Theorem 2.2 for the team semantics of the atomic formulas  $X \overset{r}{\smile}_{k\ell} y$  and  $X \overset{r}{\smile}_{k\ell} y$ .
- For the other atomic formula  $\theta$ ,  $M \models_D \theta$  iff  $M \models_s \theta$  in the usual sense for all  $s \in D$ .
- $M \models_D \phi \wedge \psi$  iff  $M \models_D \phi$  and  $M \models_D \psi$ .
- $M \models_D \exists^1 x \phi$  iff  $M \models_{D(a/x)} \phi$  for some element  $a \in M$  of data sort, where  $D(a/x) = \{s(a/x) \mid s \in D\}$ .
- $M \models_D \exists^1 r \phi$  and  $M \models_D \exists^1 \ell \phi$  are defined as above respecting the sorts of the variables.

<sup>1</sup>Notice that we assume the variable  $t$  to be always present in these atoms, although for simplicity we do not explicitly write the variable in the formulas. Cf. the team semantics given in (3).

For any set  $\Gamma \cup \{\phi\}$  of formulas, we write  $\Gamma \models \phi$  if for all models  $M$  and teams  $D$ ,  $M \models_D \gamma$  for all  $\gamma \in \Gamma$  implies  $M \models_D \phi$ . We write  $\phi \models \psi$  for  $\{\phi\} \models \psi$ .

### 3 FORMALISING MATTHEW EFFECTS

In this section, we formalise four distinct types of Matthew effects in ML and we investigate the properties of the dependence relations and Matthew effects.

*Definition 3.1 (Matthew effects).* We define the following notions:

- $y$  being subject to a **positive direct  $\ell$ -Matthew effect** with respect to  $r$  (see also Table 2(b)):

$$\text{DME}_\ell^r y ::= y \overset{r}{\mathcal{A}}_\ell y.$$

- $y$  being subject to a **positive  $x$ -mediated  $\ell$ -Matthew effect** with respect to  $r$  (see also Table 2(c)):

$$\text{MME}_\ell^r y(x) ::= x \overset{r}{\mathcal{A}}_\ell y \wedge y \overset{r}{\mathcal{A}}_\ell x.$$

- $y$  being subject to a **positive  $x$ -complete  $\ell$ -Matthew effect** with respect to  $r$  (see also Table 2(d)):

$$\text{CME}_\ell^r y(x) ::= \text{MME}_\ell^r y(x) \wedge \text{DME}_\ell^r y.$$

- $x$  and  $y$  being subjects to a **positive complete  $\ell$ -Matthew effect** with respect to  $r$  (see also Table 2(e)):

$$\text{CME}_\ell^r(x, y) ::= \text{MME}_\ell^r y(x) \wedge \text{DME}_\ell^r x \wedge \text{DME}_\ell^r y.$$

*Properties.* It is not hard to verify that the dependence relation  $X \overset{r}{\mathcal{A}}_\ell y$  satisfies the following properties:

- (Reflexivity)  $\models x \overset{r}{\mathcal{A}}_0 x$
- (Enhancing)  $x \overset{r}{\mathcal{A}}_\ell x \models \exists^1 r \ x \overset{r}{\mathcal{A}}_{k\ell} x$
- (Commutativity)  $X_1, \dots, X_n \overset{r}{\mathcal{A}}_\ell y \models X_{i_1}, \dots, X_{i_n} \overset{r}{\mathcal{A}}_\ell y$  and  $W, x_1 \otimes \dots \otimes x_n, Z \overset{r}{\mathcal{A}}_\ell y \models W, x_{i_1} \otimes \dots \otimes x_{i_n}, Z \overset{r}{\mathcal{A}}_\ell y$ , where  $i_1, \dots, i_n$  is any permutation of  $1, \dots, n$
- (Duplication)  $X_1, \dots, X_n \overset{r}{\mathcal{A}}_\ell y \models X_i, X_1, \dots, X_n \overset{r}{\mathcal{A}}_\ell y$ , where  $X_i \in \{X_1, \dots, X_n\}$
- (Projection)  $X_1, \dots, X_n \overset{r}{\mathcal{A}}_\ell y \models X_{i_1}, \dots, X_{i_k} \overset{r}{\mathcal{A}}_\ell y$ , where  $X_{i_1}, \dots, X_{i_k}$  is any subsequence of  $X_1, \dots, X_n$
- (Regrouping)  $(X \overset{r}{\mathcal{A}}_\ell y), (Z \overset{r}{\mathcal{A}}_\ell y) \models X, Z \overset{r}{\mathcal{A}}_\ell y$
- (Transitivity)  $(X \overset{r}{\mathcal{A}}_\ell y), (y \overset{r}{\mathcal{A}}_{k\ell} z) \models \exists^1 r \ X \overset{r}{\mathcal{A}}_{(k+1)\ell} z$

As a consequence of the transitivity of the dependence relation, mediated and direct Matthew effects satisfy the following properties:

- (Transitivity)  $\text{MME}_\ell^r x(y), \text{MME}_\ell^r y(z) \models \exists^1 r \ \text{MME}_{2\ell}^r x(z)$
- (Scaling)  $\text{MME}_\ell^r y(x) \models \exists^1 r_0 \ \text{DME}_{2\ell}^{r_0} x \wedge \exists^1 r_1 \ \text{DME}_{2\ell}^{r_1} y$

Moreover, a direct Matthew effect of a variable  $x$  is clearly a mediated Matthew effect where  $x$  itself is the mediator, namely,  $\text{DME}_\ell^r y \models \text{MME}_\ell^r y(y)$ , and a mediated Matthew effect is *reciprocal* for the two variables involved, namely  $\text{MME}_\ell^r y(x) \models \text{MME}_\ell^r x(y)$ .

### 4 CASE STUDY

In Section 4.1, we informally present the results of Azoulay, Stuart, and Wang (henceforth: ASW) reported in [2] about the Matthew effect in science. In Section 4.2, we formalise their analysis.

#### 4.1 Informal presentation

ASW propose an empirical test for the following proposition: scientists of higher status will have even higher status in the future. According to the definitions presented in Section 2, this is equivalent to saying that a direct Matthew effect exists with regard to a scientist's status. The empirical test revolves around the conferral to medical scientists of the prestigious title of Howard Hughes Medical Investigator (HHMI). ASW examine how the yearly number of citations for an article published by a HHMI-appointed scientist before the appointment changes as a result of the appointment. ASW assume that both the HHMI appointment and an article's yearly number of citations reflect the scientist's status.

The hypothesis that a direct Matthew effect exists with regard to a scientist's status is accepted if the articles published by HHMI appointees receive more citations after the appointment, compared to articles of similar quality published by non-appointees. The results of the analysis show that:

- The citation boost is small and it affects only the articles published up to 1 year before the HHMI appointment. Older articles do not witness a change in citations received as a result of the appointment.
- The citation boost is larger for articles published in journals with low impact factor, articles that use more novel keywords, and articles that cite a greater number of studies from other fields (i.e., that are *recombinant*). ASW argue that this is because the quality of these articles is more difficult to assess; therefore, the HHMI appointment acts as a signal of quality and more strongly affects the yearly citations these articles receive.
- The citation boost is larger for articles published by scientists who have a smaller total number of citations attached to their name or who are younger at the time of the HHMI appointment.

On the basis of these results, ASW conclude that there is a Matthew effect with regard to scientists' status, but the extent to which this is observable depends on the age of the articles published by scientists and on how easily the quality of these articles can be assessed. In addition, they conclude that the Matthew effect more strongly affects scientists who have lower status at the time they are appointed.

## 4.2 Formalisation

ASW perform a complex empirical test involving multiple variables and regression models. These are presented in detail in Appendix B. In this subsection, we first formalise ASW's statements in isolation, then we analyse the reasoning they use to draw their conclusions.

*Statements.* Each statement is presented in natural language (*text in italic*) and then formalised using ML.

ASW aim at proving that *there is a direct Matthew effect with regard to a scientist's status*. This can be formalised by the formula:

$$\exists^1 \ell \exists^1 r_a \text{ DME}_{\ell}^{r_a} \text{Status}. \quad (4)$$

ASW observe (regression  $r_1$ ) that ACitAF *positively depends on HHMI, but the effect is small* [2, page 21], which can be formalised by the formula:

$$\text{HHMI} \overset{r_1}{\text{year}} \text{ACitAF} \wedge \text{Small}(r_1, \text{ACitAF}, \{\text{HHMI}\}). \quad (5)$$

In addition, they observe (regressions  $r_2$ ,  $r_3$ , and  $r_4$ ) that *this positive dependency only affects the articles published up to 1 year before the HHMI appointment*:

$$\text{HHMI} \overset{r_2}{\text{year}} \text{ACitAF} \quad (6)$$

$$\wedge \text{Effect}(r_2, \text{ACitAF}, \{\text{HHMI}\}) \gg \text{Effect}(r_1, \text{ACitAF}, \{\text{HHMI}\}) \quad (7)$$

$$\wedge \text{Effect}(r_3, \text{ACitAF}, \{\text{HHMI}\}) \approx 0 \wedge \text{Effect}(r_4, \text{ACitAF}, \{\text{HHMI}\}) \approx 0. \quad (8)$$

Recall that the regressions  $r_2$ ,  $r_3$  and  $r_4$  are respectively based on articles published up to 1 year, 2 years, and 3 to 10 years before the HHMI appointment (see Table 4).

ASW also observe (regression  $r_5$ ) that *there is a stronger increase in citations after the HHMI appointment if the article is published in a journal with low impact factor*:

$$\text{HHMI} \otimes \text{LIF} \overset{r_5}{\text{year}} \text{ACitAF} \quad (9)$$

Furthermore, they observe that *there is a stronger increase in citations if the article is novel* (regression  $r_6$ ):

$$\text{HHMI} \otimes \text{Novel} \overset{r_6}{\text{year}} \text{ACitAF}, \quad (10)$$

or if it is recombinant (regression  $r_7$ ):

$$\text{HHMI} \otimes \text{Recombinant} \overset{r_7}{\text{year}} \text{ACitAF}. \quad (11)$$

ASW assume that *the quality of articles is more uncertain if they are published in journals with low impact factor, if they are novel, or if they are recombinant*:

$$\exists^1 \ell \exists^1 r_b \left( \{\text{LIF}, \text{Novel}, \text{Recombinant}\} \overset{r_b}{\ell} \text{UArtQ} \right). \quad (12)$$

The assumption (12) and the observations (9), (10) and (11) suggest that *the positive dependency of ACitAF on HHMI more strongly affects articles of uncertain quality*:

$$\exists^1 \ell \exists^1 r_c \left( \text{HHMI} \overset{r_c}{\ell} \text{ACitAF} \wedge \text{HHMI} \otimes \text{UArtQ} \overset{r_c}{\ell} \text{ACitAF} \right). \quad (13)$$

Moreover, ASW observe that *there is a stronger increase in the number of citations after the HHMI appointment if the article is published by a scientist who is less cited at the time of the appointment* (regression  $r_8$ ):

$$\text{HHMI} \otimes \text{Hnotwellcited} \stackrel{r_8}{\nearrow}_{\text{year}} \text{ACitAF}. \quad (14)$$

or who is younger at the time of the appointment (regression  $r_9$ ):

$$\text{HHMI} \otimes \text{Hyoung} \stackrel{r_9}{\nearrow}_{\text{year}} \text{ACitAF}. \quad (15)$$

The assumption that *the variables Hnotwellcited and Hyoung have a negative effect on Status*,

$$\exists^1 \ell \exists^1 r_d \quad (\{\text{Hnotwellcited}, \text{Hyoung}\} \searrow_{\ell}^{r_d} \text{Status}), \quad (16)$$

leads to the conclusion that *the positive dependency of ACitAF on HHMI more strongly affects scientists who have lower status at the time of the HHMI appointment than it affects scientists who have higher status*:

$$\exists^1 \ell \exists^1 r_e \quad (\text{HHMI} \stackrel{r_e}{\nearrow} \text{ACitAF} \wedge \text{HHMI} \otimes \text{Status} \searrow_{\ell}^{r_e} \text{ACitAF}). \quad (17)$$

*Reasoning.* The observations (5), (6), (7), (8) (9), (10), (11), (14) and (15) and the assumptions (12) and (16) are considered as axioms. Here we list the deduction steps ASW use to deduce (4), (13) and (17). To deduce (13), ASW use the following reasoning: *if ACitAF is subject to a positively moderated year-dependency on HHMI  $\otimes$  LIF (resp. HHMI  $\otimes$  Novel or HHMI  $\otimes$  Recombinant) witnessed by the regression  $r$ , and if UArtQ is subject to a positive LIF (resp. Novel or Recombinant) dependency, then there is an hypothetical regression  $r'$  that witnesses the fact that ACitAF is subject to a positively moderated  $\ell'$ -dependency on HHMI  $\otimes$  UArtQ.* This reasoning is captured by the following deduction step:

$$\frac{h \otimes n \stackrel{r}{\nearrow}_{\ell} c \quad \exists^1 \ell \exists^1 r \ n \ \stackrel{r}{\nearrow}_{\ell} u}{\exists^1 \ell' \exists^1 r' \ h \otimes u \ \stackrel{r'}{\nearrow}_{\ell'} c}.$$

To deduce (17), ASW use the following reasoning: *if ACitAF is subject to a positively moderated year-dependency on HHMI  $\otimes$  Hnotwellcited (resp. HHMI  $\otimes$  Hyoung) witnessed by the regression  $r_8$  (resp.  $r_9$ ), and if Hnotwellcited (resp. Hyoung) has a negative effect on Status, then there is an hypothetical regression  $r'$  that witnesses the fact that ACitAF is subject to a positively moderated  $\ell'$ -dependency on HHMI  $\otimes$  Status.* This reasoning is captured by the following deduction step:

$$\frac{h \otimes n \stackrel{r}{\nearrow}_{\ell} c \quad \exists^1 \ell \exists^1 r \ n \ \searrow_{\ell}^{r} s}{\exists^1 \ell' \exists^1 r' \ h \otimes s \ \searrow_{\ell'}^{r'} c}.$$

To deduce (4), one first need to express the assumption that HHMI is positively dependent on Status:

$$\exists^1 \ell \exists^1 r_f \quad (\text{Status} \stackrel{r_f}{\nearrow}_{\ell} \text{HHMI}). \quad (18)$$

Based on this, ASW use the following reasoning: *if HHMI is positively dependent on Status, if ACitAF is positively dependent on HHMI as witnessed by the regression  $r_1$ , then Status is positively dependent on Status, which means that Status is subject to a direct Matthew effect.* This reasoning is captured by the following deduction steps:

$$\frac{\exists^1 \ell \exists^1 r \ s \ \stackrel{r}{\nearrow}_{\ell} h \quad h \ \stackrel{r_1}{\nearrow}_{\text{year}} c}{\exists^1 \ell \exists^1 r \ s \ \stackrel{r}{\nearrow}_{\ell} c} \quad \text{and} \quad \frac{\exists^1 \ell \exists^1 r \ s \ \stackrel{r}{\nearrow}_{\ell} c \quad \exists^1 \ell \exists^1 r \ c \ \stackrel{r}{\nearrow}_{\ell} s}{\exists^1 \ell \exists^1 r \ \text{DME}_{\ell}^r s}.$$

## 5 CONCLUSION AND FURTHER RESEARCH

*Conclusion.* While there is a great deal of literature in the social sciences invoking the Matthew effect to explain important phenomena, from career dynamics to economic inequality, the concept of Matthew effect has never been properly formalised. This makes it difficult to compare and synthesise the results of different studies. This paper offers a first formalisation of the Matthew effect via a logic based on team semantics.

An original contribution of this paper is that our formalisation allows for a clear distinction between different types of Matthew effects: direct, mediated, and complete. This shows just how complicated self-reinforcing phenomena can be, because an observed Matthew effect can actually result from the interplay of direct and mediated Matthew effects. In addition, our formalisation serves a number of interrelated purposes: first, as shown in our case study, it can be used to better understand the import of empirical research and to make explicit the assumptions needed to support the researchers' conclusions; second, it can be used to compare and relate the results of different studies to one another, and thereby develop new theory on a firmer foundation; third, it allows empirical scientists to ask new research questions and formulate more precise hypotheses.



This study is only a first step in exploring logical formalisms that address the intricate phenomena concerning the Matthew effect. Further work will expand the logical analysis in the following directions:

*Studying properties of ML.* The logic ML, we introduced is defined on the basis of team semantics. A team or a data set is essentially a relation of the model, which is a second-order object. As a consequence, logics based on team semantics are usually second-order in expressive power. Indeed, the two major team-based logics, dependence logic and independence logic, are expressibly equivalent to existential second-order logic [3, 15]. The atomic dependency notions formalised in ML are more involved than the functional dependency, independence, and other dependency notions studied so far in team-based logics. Yet the language of our ML is very simple, as the only complex formulas are the conjunctions and the existentially quantified statements with weak existential quantifiers of team semantics. One natural conjecture would be that this logic is strictly weaker in expressive power than existential second-order or even first-order logic.

*A richer language with a good proof calculus.* Although we demonstrated in the case study that the simple language of ML can already express interesting facts about Matthew effects, in future work we will introduce stronger logics by expanding the language to include disjunction, negation, implication, and strong quantifiers. We also want to introduce proof calculi for these extensions or their sufficiently strong fragments.

*Reflecting the complexity of empirical tests.* Our formalisation suggests that empirical findings about Matthew effects are contingent on particular statistical analyses, performed on particular data, where particular variables are observed over particular time intervals. Choices related to research design can thus affect the empirical evidence researchers find about Matthew effects. For example, choosing to observe the variables yearly rather than monthly, weekly, or daily can determine whether a direct Matthew effect is found in the place of a mediated one. Moreover, the fact that some variables like quality and uncertainty about quality remain unobserved can conceal important dependencies; as a result, a Matthew effect may appear to be mediated by a certain (observed) variable whereas in fact it is mediated by another (unobserved) one, which depends or is dependent on the apparent mediator. In formalising Matthew effects, or indeed any dependency tested via statistical analysis, one must be able to express these details within the syntax of the logic.

## REFERENCES

- [1] C Antonelli and F Crespi. 2013. The ‘Matthew effect’ in R&D public subsidies: The Italian evidence. *Technological Forecasting and Social Change* 80, 8 (2013), 1523–1534.
- [2] P Azoulay, T Stuart, and Y Wang. 2013. Matthew: Effect or fable? *Management Science* 60, 1 (2013), 92–109.
- [3] Pietro Galliani. 2012. Inclusion and Exclusion in Team Semantics: On some logics of imperfect information. *Annals of Pure and Applied Logic* 163, 1 (January 2012), 68–84.
- [4] Pietro Galliani and Lauri Hella. 2013. Inclusion Logic and Fixed Point Logic. In *Computer Science Logic 2013 (CSL 2013)*, CSL 2013, September 2–5, 2013, Torino, Italy. 281–295. DOI : <https://doi.org/10.4230/LIPIcs.CSL.2013.281>
- [5] Erich Grädel and Jouko Väänänen. 2013. Dependence and Independence. *Studia Logica* 101, 2 (April 2013), 399–410.
- [6] J. Hamilton. 1994. *Time Series Analysis*. Princeton University Press.
- [7] Miika Hannula. 2015. Axiomatizing first-order consequences in independence logic. *Annals of Pure and Applied Logic* 166, 1 (2015), 61–91.
- [8] W. Hodges. 1997. Compositional Semantics for a Language of Imperfect Information. *Logic Journal of the IGPL* 5 (1997), 539–563.
- [9] W. Hodges. 1997. Some Strange Quantifiers. In *Structures in Logic and Computer Science: A Selection of Essays in Honor of A. Ehrenfeucht*, J. Mycielski, G. Rozenberg, and A. Salomaa (Eds.). Lecture Notes in Computer Science, Vol. 1261. London: Springer, 51–65.
- [10] J W Kim and B G King. 2014. Seeing stars: Matthew effects and status bias in Major League Baseball umpiring. *Management Science* 60, 11 (2014), 2619–2644.
- [11] Juha Kontinen and Jouko Väänänen. 2013. Axiomatizing first-order consequences in dependence logic. *Annals of Pure and Applied Logic* 164, 11 (2013).
- [12] R K Merton. 1968. The Matthew effect in science. *Science* 159, 3810 (1968), 56–63.
- [13] R K Merton. 1988. The Matthew effect in science, II: Cumulative advantage and the symbolism of Intellectual Property. *Isis* 79, 4 (1988), 606–623.
- [14] D Rigney. 2010. *The Matthew effect: How advantage begets further advantage*. Columbia University Press, New York City, NY.
- [15] Jouko Väänänen. 2007. *Dependence Logic: A New Approach to Independence Friendly Logic*. Cambridge: Cambridge University Press.
- [16] Fan Yang. 2016. Negation and Partial Axiomatizations of dependence and independence logic revisited. In *Proceedings of 23rd Workshop on Logic, Language, Information and Computation (LNCS 9803)*, Jouko Väänänen, Asa Hirvonen, and Ruy de Queiroz (Eds.). Springer-Verlag, 410–431.

## A APPENDIX - MATTHEW EFFECTS

### A.1 Toy example of a Matthew effect

Matthew effects are often detected by researchers while analyzing empirical data. In the statistical literature, these represent a form of autocorrelation [6]. Table 1 presents a hypothetical dataset that shows *prima facie* evidence of a Matthew effect. This data concerns the careers of visual artists: each row contains information about an artist during a given year. The first column includes the artist ID; the second column includes the number of artworks sold by the artist during the observation year; the third column includes the number of times the artist or her work were reviewed by the media during the observation year; the fourth column specifies the observation year. Assume that all artists started their career in 2010 and that they were equally productive during the study period. The data suggests that Artist A started accumulating reviews from the very beginning, and sales quickly followed. A similar pattern can be observed for Artist B, though with some delay. In the case of Artist C, however, this trend never began.

We may presume that both sales and reviews are subject to a Matthew effect because selling more artworks leads to greater odds of selling artwork in the future. Similarly, being reviewed increases the odds of future reviews. However, it is also possible—and indeed highly likely—that being reviewed increases the odds of future sales, and that selling artwork increase the odds of future reviews. The fact that these dependencies occur at the same time makes the individual effects difficult to isolate.

Table 1. A data set (or a team)

Artist	Sales	Reviews	Time	...	Artist	Sales	Reviews	Time	...	Artist	Sales	Reviews	Time	...
A	0	1	2010	...	B	0	0	2010	...	C	0	0	2010	...
A	1	2	2011	...	B	0	0	2011	...	C	1	1	2011	...
A	1	1	2012	...	B	0	1	2012	...	C	0	0	2012	...
A	2	4	2013	...	B	0	2	2013	...	C	0	1	2013	...
A	4	7	2014	...	B	2	5	2014	...	C	1	2	2014	...
A	7	9	2015	...	B	4	8	2015	...	C	0	1	2015	...
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

### A.2 Different Matthew effects

Table 2 contains a graphic representation of the definitions of the positive dependency introduced in Theorem 2.1 and of the Matthew effects defined in Section 3. To help understand the definitions of the different Matthew effects and Table 2, we present a simple example of positive dependency.

*Example A.1.* Consider a very simple dependence relation  $x \nearrow_{\ell}^r y$  with a linear regression function

$$y_{(t)} = \alpha + \beta x_{(t-\delta)} + \gamma_1 w_{1(t-\delta)} + \cdots + \gamma_l w_{l(t-\delta)} + \epsilon,$$

where  $\beta$  is significantly greater than 0. Since we also have

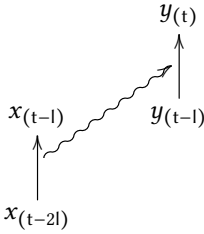
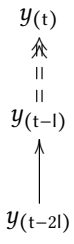
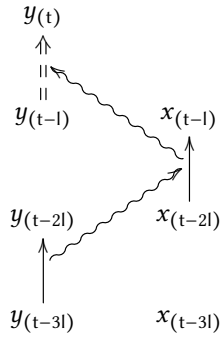
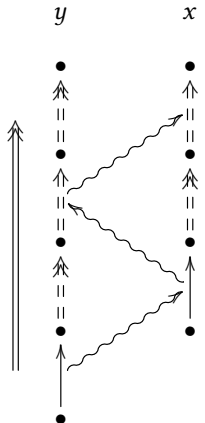
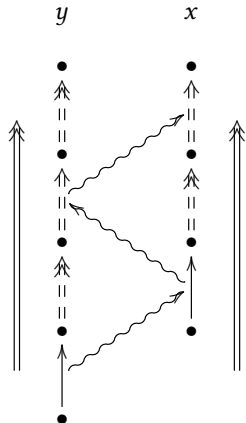
$$y_{(t-\delta)} = \alpha + \beta x_{(t-2\delta)} + \gamma_1 w_{1(t-2\delta)} + \cdots + \gamma_l w_{l(t-2\delta)} + \epsilon,$$

when all the independent variables except  $x$  are held constant, we obtain

$$y_{(t)} - y_{(t-\delta)} = \beta(x_{(t-\delta)} - x_{(t-2\delta)}),$$

meaning that the value of  $y$  increases as the value of  $x$  increases (see also Table 2(a)).

Table 2. Types of Matthew effect

(a) Dependence relation	
	
(b) Direct Matthew effect	(c) Mediated Matthew effect
	
(d) Complete Matthew effect for $y$	(e) Complete Matthew effect for $(x, y)$
	

## B APPENDIX - CASE STUDY

### B.1 Variables

The key variables used by ASW to analyze the Matthew effect among scientists in [2] is reported in Table 3. Some of these variables measure relevant characteristics of a scientist, including Author, HHMI, Hdate, Hage, and Hcit. Other variables measure characteristics related to the article published by the scientist, including Article, ArtY, Journal, NbAut, Apos, ACitAF, ACitBF, IF, Novelty, and Recombination.

In [2], the variables Hage, Hcit, IF, Novelty, and Recombination are split along the median value observed in the data, giving rise to the following boolean variables: Hyoung, which is true if the scientist is younger than the median at the time of the HHMI appointment; Hwellcited, which is true if the scientist has a greater number of total citations than the median at the time of the HHMI appointment; HIF, which is true if the journal where the article is published has higher impact factor than the median; Novel, which is true if the keywords associated with the article are more novel than the median; and Recombinant, which is true if the proportion of out-of-field literature cited by the article is greater than the median.

In addition to these variables, which are measured and actually used in ASW’s statistical analysis, the claims made in [2] about the Matthew effect involve a number of additional variables, which cannot be measured directly but are assumed to depend on some of the observed variables. These include Status, i.e., the status of the focal scientist, as well as ArtQ and UArtQ, which represent the quality of the scientist’s article and the uncertainty about the quality of the scientist’s article, respectively.

### B.2 Regressions

A regression is performed via the application of an algorithm, i.e., an *estimator*, to the observed data. The algorithm yields a set of *coefficients*, which correspond to the  $\beta$  mentioned in Section 2. Each coefficient represents an *effect*, i.e., a change in the value of the dependent variable  $y$  that results from a one-unit increase in the value of the independent variable  $x$ . Each coefficient is associated with a *level of statistical significance*, i.e., a value between 0 and 1 that represents the probability of observing the estimated effect in the data. The lower this value, the greater the probability of observing the effect. A coefficient with a level of statistical significance below some predetermined threshold is said to be *significantly greater than 0* if the coefficient is positive, and *significantly smaller than 0* if the coefficient is negative. If the level of statistical significance is above the predetermined threshold, the coefficient is said to be *non-significant* or *equivalent to zero*. In [2], the chosen significance threshold is 0.05, as is conventional in the social sciences.

Table 4 reports the full list of the regressions performed by ASW in [2]. These are indexed by numbers 1–9. These regressions concern the observed variables listed in Table 3. In addition, Table 4 reports a list regressions that are not actually performed by ASW in [2], but their results are nonetheless relevant to ASW’s analysis. We call these *hypothetical regressions*. These are not actually performed because they concern the unobserved variables listed in Table 3. However, they could be performed if it were possible to observe these variables. These hypothetical regressions are indexed by letters *a–e*. In every regression, the dependent variable is ACitAF.

Table 3. List of variables

		Variable	Meaning	Type	Time dependent
<b>Author related</b>	<b>Observed</b>	Author	Scientist ID	$\mathbb{N}$	
		HHMI	Author is appointed HHMI	0/1	
		Hdate	Date of the HHMI appointment	$\mathbb{N}$	
		Hage	Age of the author at HHMI appointment	$\mathbb{N}$	
		Hyoung	Author was young at HHMI appointment	0/1	
		Hold	Author was old at HHMI appointment	0/1	
		Hcit	Author's total citations at HHMI appointment	$\mathbb{N}$	
		Hwellcited	Author was well cited at HHMI appointment	0/1	
		Hnotwellcited	Author was not well cited at HHMI appointment	0/1	
	<b>Unobserved</b>	Status	Status of the author	$\mathbb{R}$	
<b>Article related</b>	<b>Observed</b>	Article	Article ID	$\mathbb{N}$	
		ArtY	Date of article publication	$\mathbb{N}$	
		Journal	Journal ID	$\mathbb{N}$	
		NbAut	Number of authors	$\mathbb{N}$	
		Apos	Position of the focal author in the list of article authors	$\mathbb{N}$	
		ACitBF	Yearly number of article citations at author's HHMI appointment	$\mathbb{N}$	yes
		ACitAF	Yearly number of article citations after author's HHMI appointment	$\mathbb{N}$	yes
		IF	Impact factor of the journal	$\mathbb{R}$	
		HIF	Journal has high impact factor	0/1	
		LIF	Journal has low impact factor	0/1	
		Novelty	Novelty of the article (i.e., mean age of the keywords)	$\mathbb{R}$	
		Novel	Article is novel	0/1	
		NotNovel	Article is not novel	0/1	
		Recombination	Level of recombination of the article (i.e., proportion of out-of-field literature cited)	$\mathbb{R}$	
		Recombinant	Article is recombinant	0/1	
	NotRecombinant	Article is not recombinant	0/1		
	<b>Unobserved</b>	ArtQ	Quality of the article	$\mathbb{R}$	
		UArtQ	Uncertainty about the quality of the article	$\mathbb{R}$	

Table 4. List of regressions

<b>Regression</b>	<b>Description</b>
$r_1$	Based on the full sample. Independent variables include all the observed variables listed in Table 3
$r_2$	Like $r_1$ , but the sample includes only articles published up to 1 year before the HHMI appointment
$r_3$	Like $r_1$ , but the sample includes only articles published 2 years before the HHMI appointment
$r_4$	Like $r_1$ , but the sample includes only articles published 3 to 10 years before the HHMI appointment
$r_5$	Like $r_2$ , but estimating a different effect of HHMI when HIF is true or false
$r_6$	Like $r_2$ , but estimating a different effect of HHMI when Novel is true or false
$r_7$	Like $r_2$ , but estimating a different effect of HHMI when Recombinant is true or false
$r_8$	Like $r_2$ , but estimating a different effect of HHMI when Hwellcited is true or false
$r_9$	Like $r_2$ , but estimating a different effect of HHMI when Hyoung is true or false
$r_a$	Hypothetical regression where the direct Matthew effect on Status is estimated
$r_b$	Hypothetical regression where the effects of HIF, Novel, and Recombinant on UArtQ are estimated
$r_c$	Hypothetical regression where the effect of UArtQ on ACitAF is estimated
$r_d$	Hypothetical regression where the effects of Hwellcited and Hyoung on Status are estimated
$r_e$	Hypothetical regression where the effect of Status on ACitAF is estimated
$r_f$	Hypothetical regression where the effect of Status on HHMI is estimated