

# The dynamics of higher-order novelties

Gabriele Di Bona<sup>1, 2, 3, 4</sup>, Alessandro Bellina<sup>3, 4, 5</sup>, Giordano De Marzo<sup>4, 5, 6, 7</sup>, Angelo Petralia<sup>8</sup>, Iacopo Iacopini<sup>9, 10</sup>, and Vito Latora<sup>1, 7, 11, \*</sup>

<sup>1</sup>School of Mathematical Sciences, Queen Mary University of London, London E1 4NS, United Kingdom

<sup>2</sup>CNRS, GEMASS, 59 rue Pouchet, F-75017, Paris, France

<sup>3</sup>Sony Computer Science Laboratories Rome, I-00184, Rome, Italy

<sup>4</sup>Centro Ricerche Enrico Fermi, I-00184 Rome, Italy

<sup>5</sup>Dipartimento di Fisica Università “Sapienza”, I-00185 Rome, Italy.

<sup>6</sup>Sapienza School for Advanced Studies, “Sapienza”, I-00185 Rome, Italy.

<sup>7</sup>Complexity Science Hub Vienna, A-1080 Vienna, Austria

<sup>8</sup>Department of Economics and Business, University of Catania, I-95128 Catania, Italy

<sup>9</sup>Network Science Institute, Northeastern University London, London, E1W 1LP, United Kingdom

<sup>10</sup>Department of Physics, Northeastern University, Boston, MA 02115, USA

<sup>11</sup>Dipartimento di Fisica ed Astronomia, Università di Catania and INFN, I-95123 Catania, Italy

\*Email: v.latora@qmul.ac.uk

## ABSTRACT

Studying how we explore the world in search of novelties is key to understand the mechanisms that can lead to new discoveries. Previous studies analyzed novelties in various exploration processes, defining them as the first appearance of an element. However, novelties can also be generated by combining what is already known. We hence define higher-order novelties as the first time two or more elements appear together, and we introduce higher-order Heaps’ exponents as a way to characterize their pace of discovery. Through extensive analysis of real-world data, we find that processes with the same pace of discovery, as measured by the standard Heaps’ exponent, can instead differ at higher orders. We then propose to model an exploration process as a random walk on a network in which the possible connections between elements evolve in time. The model reproduces the empirical properties of higher-order novelties, revealing how the network we explore changes over time along with the exploration process.

## Introduction

As humans, we experience novelties as part of our daily life. By the term *novelty* we generally indicate two apparently different things<sup>1</sup>. On the one hand, we can think of a novelty as the first time we visit a neighborhood, enter a newly launched pub, or listen to a song from an artist we previously did not know. In this case, the novelty represents a discovery for a single individual of a place, an artist or, more in general, an item. On the other hand, there are discoveries that are new to the entire population, as could be a technological advancement or the development of a new drug. However, these two cases are not entirely distinct, as the second set of novelties, those new to everyone, represent a subset of the first one. Analysing how novelties emerge, both at the individual level and at the level of the entire population, is key to understand human creativity and the neural and social mechanisms that can lead to new discoveries.

The increasing availability of data on human behavior and consumption habits has allowed to study how humans explore the world, how novelties emerge in different contexts, and how they are distributed in time<sup>1-3</sup>. Empirical investigations cover a broad range of different areas<sup>4</sup>, ranging from science<sup>5</sup> and language<sup>6,7</sup>, to gastronomy<sup>8</sup>, goods or products<sup>9</sup>, network science<sup>10</sup>, information<sup>11</sup>, and cinema<sup>12</sup>, to name a few relevant examples. No matter the topic, one can always

represent data coming from real-world exploration processes as sequences of elements or “items” that are sequentially adopted or consumed<sup>13</sup>. For instance, the activity of a user on an online digital music platform is turned into a sequence of listened songs, and a novelty is defined as the first time a song, or an artist, appears in the sequence<sup>14</sup>. Analogously, articles published in a scientific journal can be turned into a time-ordered sequence of concepts or keywords discovered by the community, and a novelty can be defined, again, as the first-time appearance of a keyword<sup>3</sup>. Under this framework, evidence shows that—independently of the system they belong to—novelties seem to obey the same statistical patterns in the way they are distributed and correlated in time<sup>1</sup>. Indeed, a long tradition of works, started by the Yule-Simon processes for text generation<sup>15,16</sup>, shows that most empirical sequences follow Heaps’<sup>17-19</sup>, Zipf’s<sup>20-24</sup>, and Taylor’s laws<sup>25</sup>.

Along with data-driven investigations, a relevant scientific problem is that of finding plausible mechanisms to reproduce and explain the empirical observations. What are the drivers controlling the appearance of new items in a sequence? How do humans explore the seemingly infinite space of possibilities in search of novelties? Interestingly, an insightful answer comes from biology, where, in 1996, Stuart Kauffman introduced the concept of the *adjacent possible*<sup>26</sup> (AP) referring to “*all those molecular species that are not members of the actual, but are one reaction step away from the actual*”. In-

spired by previous works by Packard and Langton<sup>27–29</sup>, the AP provides a fresh view on the problem, for which discoveries (the possible) can only be found among those items which are close (the adjacent) to what is already known (the actual). New discoveries would then generate an expanding space of opportunities that are only available to us in the moment we “unlock” what is adjacent to them. Kauffman’s AP has seen many interesting applications ranging from biology<sup>26,30</sup> and economics<sup>9,31</sup> to models of discovery and innovation<sup>1,3,6</sup>. Among these, of particular interest is the recently proposed Urn Model with Triggering (UMT)<sup>1,6,32</sup>. Building upon the work of Pólya<sup>33,34</sup>, the UMT adds to the traditional *reinforcement* mechanism of the Pólya urn’s scheme a *triggering* mechanism that expands the space of possible discoveries upon the extraction of each novelty. Being able to reproduce the empirical laws and thanks to its simplicity, the UMT has been used to study various systems with an expanding set of “items”, like the rise and fall of popularity in technological and artistic productions<sup>2</sup>, the emergence and evolution of social networks<sup>35</sup>, and the evolution of the cryptocurrency ecosystem<sup>36</sup>. One could also picture ideas, concepts, or items as the linked elements of an abstract network. In this view, the exploration process can be modelled as a random walk over this network, where the AP accounts for the emergence of the new starting from the “edge of what is known” within the network. Approaches based on random walks have been used to investigate the cognitive growth of knowledge in scientific disciplines<sup>3</sup>, and further extended to account for multi-agent systems, where the individual exploration of the agent is enriched by social interactions<sup>14,37</sup>.

The idea of the AP, modelled either in terms of extractions from urns or random walks over a network, is of great importance to understand the processes leading to novelties. There is, however, another important mechanism of creation of the new which is neglected by the frameworks discussed above: novelties can arise from the combination of already-known elements. For instance, a meaningless sequence of words, if ordered in a different way, may generate elegant poetry<sup>38,39</sup>. Novel combinations of existing hashtags may lead to new social-media trends<sup>40,41</sup>. Different orderings of the same musical notes may in principle generate an endless number of songs<sup>42</sup>. The mechanics of combination of “pre-existing” items has been studied in various fields, e.g., in biology where new associations of various entities produce new organisms. It has been shown that the immune system recombines existing segments of genes to produce new receptors<sup>43,44</sup>. Also, publications and collaborations in science<sup>45</sup> are typically combinations of research ideas<sup>46–48</sup> and expertise<sup>49–51</sup>. Similarly, in innovation economics, as originally discussed by Schumpeter<sup>52,53</sup> and confirmed by recent works on the generation of technologies<sup>54–56</sup>, new combinations of existing factors, that interact in a technological production process, may give rise to innovations, which rule out of the market obsolete products and services<sup>57,58</sup>, thus increasing the probability of reaching further innovations (the so-called

“creative destruction”).

In this context, the aim of this paper is to explore a more general notion of novelty, including novel combinations of existing elements. We thus investigate the dynamics of “higher-order” novelties, i.e., novel pairs, triplets, etc., of items in a sequence. In particular, we focus on the Heaps’ law, which characterizes the growth of the number of novelties in the sequence as a power-law, whose exponent is a proxy for the rate of discovery<sup>18</sup> in the related process. Namely, we introduce higher-order Heaps’ laws to characterize the rate at which novel combinations of two and more elements appear in a sequence. We then analyse various types of empirical sequences, ranging from music listening records, to words in texts, and concepts in scientific articles, finding that Heaps’ laws also hold at higher orders. We discover that processes with a similar rate of discovery of single items can instead display different rates of discovery at higher orders, and can hence be differentiated by looking at higher-order novelties. We therefore propose a new model which is capable of reproducing various empirically observed features of higher-order Heaps’ laws. In our model the process of exploration is described as an edge-reinforced random walk with triggering (ERRWT) on a network. In our framework, the novelties at different orders (nodes and links visited for the first time by the walker) shape the growth of the network by reinforcing traversed links, while triggering the addition of new elements through the expansion and exploration of the adjacent possible. This expansion can happen whenever a node is visited for the first time, making other nodes accessible to the explorer, but also whenever a link is firstly used. In this case, the newly established connection will trigger novel combinations between previously explored nodes. By fitting the contributions of the two mechanisms of reinforcement and triggering, the ERRWT model is able to reproduce well the variety of scaling exponents found in real systems for the Heaps’ laws at different orders.

## Results

### Higher-order Heaps’ laws

An exploration process can be represented as an ordered set of  $T$  symbols  $\mathcal{S} = \{a_1, a_2, \dots, a_T\}$  sequentially explored. Such a set describes the sequence of “events” or “items” produced along the journey, e.g., the songs listened by a given individual over time, the list of hashtags posted on an online social network, the list of words in a text, or any other ordered list of items or ideas generated by single individuals or social groups<sup>1,13,37,59</sup>. Similarly, in the context of some recent modelling schemes of discovery, the sequence  $\mathcal{S}$  can be made of the colors of balls extracted from an urn<sup>1,37,59</sup>, or the nodes visited over time by a random walker moving on a network<sup>3</sup>. Although real-world events have an associated time, here, for simplicity, we focus only on their sequence, i.e., the relative temporal order of the events, neglecting the precise time at which they happen. For instance, if a person listens to song  $a_1$  at time  $t_1$ , song  $a_2$  at time  $t_2$ , song  $a_i$  at time  $t_i$ , and so on,

with  $t_1 < t_2 < \dots < t_i < \dots$ , we neglect these times and only retain the order of the songs in the sequence  $\{a_1, a_2, \dots, a_T\}$ . In other words, we assume that  $a_1$  is associated to the discrete time  $t = 1$ ,  $a_2$  is associated to time  $t = 2$ , and so forth.

Among the different ways to characterize the discovery rate of a given process, the Heaps' law,  $D(t) \sim t^\beta$ , describes the power-law growth of the number of novelties as a function of the number of items in the sequence, i.e., how the number  $D(t)$  of novel elements in the sequence  $\mathcal{S}$  scale with the sequence length  $t$ <sup>18</sup>. The so-called (standard) Heaps' exponent  $\beta$ , that from now on we indicate as *1<sup>st</sup>-order Heaps' exponent*  $\beta_1$ , is thus a measure of the pace of discovery of the process that generated the considered sequence. Given that the number of different elements  $D_1(t) \equiv D(t)$  is smaller (or equal) than the total length  $t$  of the sequence, the value of  $\beta_1$  is always bounded in the interval  $[0, 1]$ , with the extreme case  $\beta_1 = 1$  reached by a process that generates new elements at a linear rate.

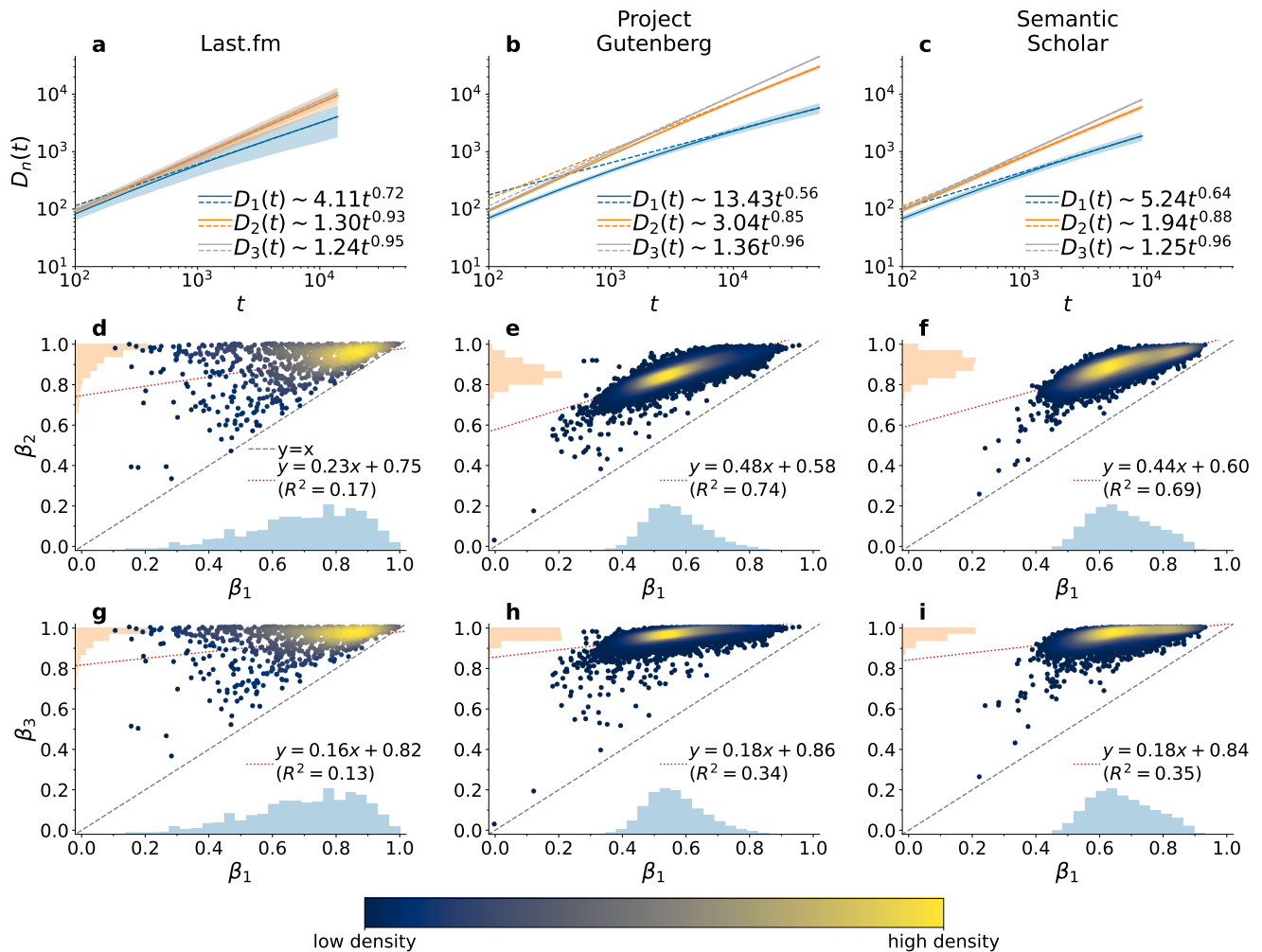
Here, we propose to go one step beyond and look at novelties as novel pairs, triplets, and higher-order combinations of consecutive symbols in a sequence<sup>60</sup>. For instance, when exploring a network, a novel pair is represented by the first visit of a link. In order to measure the pace of discovery of these higher-order compounds starting from a sequence of events  $\mathcal{S}_1 \equiv \mathcal{S}$ , we first create the surrogate sequence of overlapping pairs  $\mathcal{S}_2 = \{(a_1, a_2), (a_2, a_3), \dots, (a_{T-1}, a_T)\}$ . Considering for example the sentence “*One ring to rule them all*”, from the sequence of events  $\mathcal{S}_1 = \{one, ring, to, rule, them, all\}$  we obtain the sequence of overlapping pairs  $\mathcal{S}_2 = \{(one, ring), (ring, to), (to, rule), (rule, them), (them, all)\}$ . From  $\mathcal{S}_2$  we can then compute the number  $D_2(t)$  of different pairs among the first  $t$  ones, with  $t \leq T - 1$ . Notice that, in this manuscript, we consider the pairs *(one, ring)* and *(ring, one)* as two different pairs, i.e., order matters. By construction, we always have  $D_1(t) \leq D_2(t) \leq t$ , since, on the one hand, for each new element added to  $\mathcal{S}_1$  there is a new pair in  $\mathcal{S}_2$ , and, on the other hand, there cannot be more than  $t$  different pairs among  $t$  items. From the power-law scaling  $D_2(t) \sim t^{\beta_2}$ , we can then extract the value of  $\beta_2$ , which we refer to as the *2<sup>nd</sup>-order Heaps' exponent*. This definition can be naturally extended to any order  $n$ , considering the sequence  $\mathcal{S}_n$  of consecutive overlapping  $n$ -tuples present in  $\mathcal{S}_1$ . Notice that, if  $|\mathcal{S}_1| = T$ , then  $|\mathcal{S}_n| = T - n + 1$ . We can hence compute the number  $D_n(t)$  of different tuples among the first  $t$  tuples in  $\mathcal{S}_n$ , and extract the *n<sup>th</sup>-order Heaps' exponent*  $\beta_n \in [0, 1]$  from  $D_n(t) \sim t^{\beta_n}$ . Notice also that the *n<sup>th</sup>-order Heaps' exponent* can also be interpreted as the first order Heaps' exponent of a sequence whose events are the overlapping  $n$ -tuples of the original sequence. Finally, it is worth remarking that such an approach is close to the analysis of Zipf's law in linguistic data for  $n$ -grams or sentences<sup>61,62</sup>. In this context, studies showed that as one moves from graphemes, to words, sentences, and  $n$ -grams, the Zipf's exponent (reciprocal of the Heaps' exponent for infinitely long sequences<sup>19</sup>) gradually diminishes. This implies that  $n$ -grams or sentences are characterized by a

larger novelty rate than words, a behavior analogous to what we have discussed above.

### Analysis of real-world data sequences

We start investigating the emergence of novelties of different orders in empirical exploration processes associated to three different data sets. These data sets are substantially different in nature, since they refer, respectively, to songs listened by users of *Last.fm*, words in books collected in the *Project Gutenberg*, and words of titles of scientific journals from *Semantic Scholar* (more details on the data can be found in *Materials and Methods*). In Fig. 1(a-c) we plot the average temporal evolution of the number  $D_n(t)$  of novelties of order  $n$ , with  $n = 1, 2, 3$ , in the three data sets (from left to right, respectively, Last.fm, Project Gutenberg, Semantic Scholar). In order to avoid spurious effects due to different lengths of the sequences, we restrict these averages to the sequences of length  $T$  greater than the median length  $\tilde{T}$  in the corresponding data set (see Fig. S1 in the Supplementary Information (SI) for their distribution). Each continuous curve, plotted up to length  $\tilde{T}$ , is obtained by averaging  $D_n(t)$  over all such sequences, while the shaded area represents one standard deviation above and below the mean. We also perform power-law fits (see *Materials and Methods* for details on the procedure), and plot the resulting curves as dashed lines, with the fitted function shown in the legend. Focusing first on the broadly-studied (1<sup>st</sup>-order) Heaps' law, notice how the power-law fit is only accurate in the last part of the sequence. This highlights that the Heaps' law starts after a transient phase, where most of the events are new for the individual, as also reported in Ref.<sup>1</sup> and similarly reported in other contexts<sup>63-67</sup>. Secondly, notice how the *n<sup>th</sup>-order Heaps' law*, with  $n = 2, 3$ , is valid across the data sets, but with different values of the fitted exponents, especially for  $n = 2$ . Finally, as expected from their definition, the fitted Heaps' exponents of order  $n + 1$ , i.e.,  $\beta_{n+1}$ , are higher than the lower-order ones, that is,  $\beta_{n+1} \geq \beta_n$ .

To explore the gain in information brought by the higher-order Heaps' exponents with respect to the 1<sup>st</sup>-order Heaps', we now look directly at individual sequences. In Figure 1(d-i) we show the scatter plots of  $\beta_2$  (d-f) and  $\beta_3$  (g-i) against  $\beta_1$ , where each point refers to a single sequence from Last.fm (d,g), Project Gutenberg (e,h), or Semantic Scholar (f,i), with colors representing how dense points are (see color bar at the bottom of the figure). Here, we filter out sequences whose fitted exponent has a standard error above the 0.05 threshold (see Table S1 in SI for more details), for which the Heaps' law cannot be considered valid. This filtering removes only 30 (3.37%), 8 (0.04%), and 5 (0.03%) sequences in the three data sets, respectively. Furthermore, we have removed sequences for which the extracted value of  $\beta_2$  is higher than the associated value of  $\beta_1$ , or for which  $\beta_3 > \beta_2$ , since  $D_n(t) \leq D_{n+1}(t)$  as previously discussed. This filtering removes 53 (6.16%), 7 (0.04%), 6 (0.03%) in the three data sets, retaining a total of 807, 19 622, and 18 909 sequences, respectively. Looking at Figure 1(d), we see how users of Last.fm sharing the same



**Figure 1. Higher-order Heaps' exponents in real-world data sets.** (a-c) Average number  $D_n(t)$  of novelties of order  $n$ , with  $n = 1, 2, 3$ , as a function of the sequence length  $t$ , and fit of the associated Heaps' laws (dashed lines), with estimated exponents shown in the legend. Shaded area represents one standard deviation above and below the average. (d-i) Scatter plots between the (1<sup>st</sup>-order) Heaps' exponents  $\beta_1$  and the  $n^{\text{th}}$ -order exponents  $\beta_n$ , with  $n = 2$  (d-f) and 3 (g-i). Each point refers to a different sequence, with colors representing the density of points (see color bar). Each panel also reports histograms of exponents distributions, the bisector  $y = x$  (dashed gray line), as well as the fitted linear model (dotted red line) with the value of its coefficient of determination  $R^2$ . Each column refers to a different data set: (a,d,g) Last.fm, (b,e,h) Project Gutenberg and (c,f,i) Semantic Scholar, respectively.

value of  $\beta_1$  can have very different values of  $\beta_2$ . Conversely, the other two data sets present stronger correlation between  $\beta_2$  and  $\beta_1$ . To quantitatively characterize this, we fit a linear model with an ordinary least squares method, displayed in each plot as a red dotted line. In the legend we also report the value of the related coefficient of determination  $R^2$ , which represents the percentage of variance of the dependent variable explained by the linear fit with the independent variable. For the Last.fm data set, points are much more spread around the linear fit compared to the other two data sets, as also confirmed by the values of  $R^2$ , indicating a greater variability in listening habits compared to writing habits. Moreover, the values of the parameters of the linear fit greatly change across data sets and orders. In particular, Last.fm is characterized by a much lower slope and intercept compared to the other data sets for the same order. Furthermore, we notice how, for each data set, the values of  $\beta_3$  are much higher and more spread than the respective values of  $\beta_1$  and  $\beta_2$ , resulting in lower values of  $R^2$  in the linear fit between  $\beta_1$  and  $\beta_3$ .

At an aggregate level, we observe that at all orders the distribution of the Heaps' exponents are very different across data sets (see Fig. S2 in SI for a comparative figure, while further statistical information on the Heaps' exponents distribution can be found in Table S2 in SI). The exponents are more spread in Last.fm, which also shows a higher average of  $\beta_1$  and  $\beta_2$ , but a lower one for  $\beta_3$  compared to the other data sets. Distributions for Project Gutenberg and Semantic Scholar, which are both related to linguistic data, are more peaked. Such peaks appear at higher values for the latter data set. This could be the result of how titles of scientific papers are written with respect to books or poems, that is, concentrating the whole message of a scientific work in a few significant and specialised words, avoiding stop-words and repetition. In addition, scientific advancements tend to favor the combinations of previously existing scientific concepts to form new ones, while the same does not apply to non-scientific literature in general, where instead similar constructions tend to be repeated across the piece.

Finally, similar results are obtained also for more coarse-grained sequences generated by using artists and stemmed words instead of songs and words (see Fig S3 in SI). Furthermore, in Fig S4 in SI we analyse higher-order novelties in the collective sequences obtained by randomly concatenating all the individual sequences of each data set<sup>1</sup>.

### Analysis of existing models

After studying higher-order Heaps' laws in real data, we check whether the observed patterns can be reproduced by the available models of discovery processes. We start this analysis from the Urn Model with Triggering (UMT). In such a model, sequences of events are generated by the extraction of colored balls from an urn<sup>1</sup>, where different colors correspond to different events or items being discovered or adopted. Here, an event in the sequence is simply represented by the color extracted. In the UMT, for each extracted ball, the correspond-

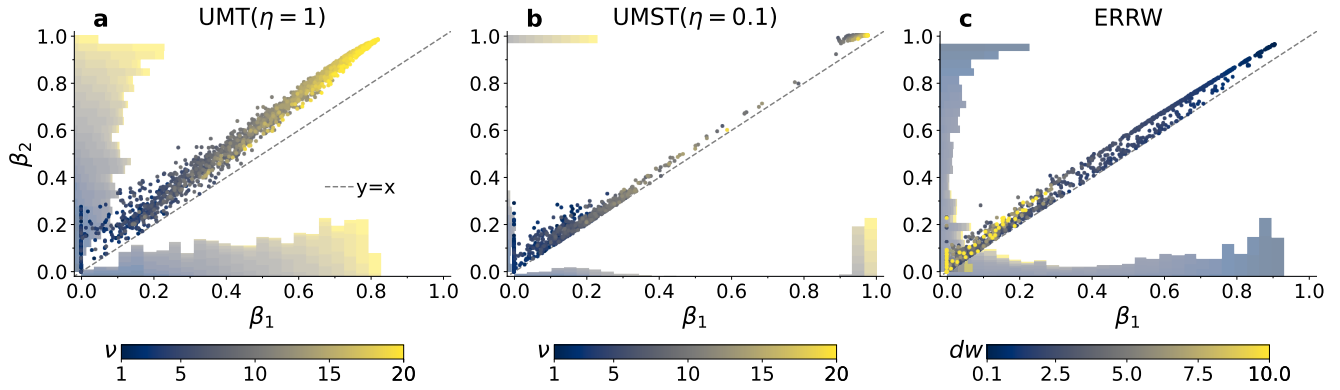
ing color is reinforced by adding  $\rho$  additional balls, of the same color, to the urn. At the same time, whenever a novel color is drawn, the discovery triggers the addition of  $v + 1$  balls of new different colors to the urn (see detailed model definition in *Materials and Methods*). The reinforcement process ensures the wide-spread adoption of items or concept that were frequently adopted in the past. Conversely, the triggering mechanism mimics the adjacent possible expansion, since each novelty makes the space of possible colors expand. Intuitively, these two parameters modulate the exploit-explore tendency of the system, with a more pronounced exploratory behavior for larger  $v/\rho$  ratios.

Previous studies have shown that the 1<sup>st</sup>-order Heaps' law is verified in sequences generated by UMT simulations<sup>1,6</sup>. In particular, the number of novelties in the model grows asymptotically as  $D_1(t) \sim t^{\frac{v}{\rho}}$  when  $v < \rho$ , while a linear behavior is found for  $v > \rho$ . We hence focus on the most interesting case  $v \leq \rho$ , studying how variations of the two parameters  $\rho$  and  $v$ , respectively representing the reinforcement and the increase in size of the adjacent possible, affect the Heaps' law at various orders. Since the pace of discovery effectively depends only on the fraction  $v/\rho$ , we fix  $\rho = 20$  and numerically simulate the UMT with  $v = 1, 2, 3, \dots, 20$  for  $T = 10^5$  time-steps, obtaining sequences of length comparable to the data sets (see Fig. S1 in SI). For each set of parameters we run 100 simulations, generating a total of  $2 \times 10^3$  synthetic sequences. Then, for each generated sequence, we compute the temporal evolution of the number of novelties  $D_n(t)$ , and estimate a power-law fit, extracting the related  $n^{\text{th}}$ -order Heaps' exponent  $\beta_n$ . In Fig. 2(a), we show how the extracted values of  $\beta_2$  change with respect to  $\beta_1$  across simulations. The color represents the value of the parameter  $v$ , as indicated by the color bar. We observe that, although the exponents span the interval  $(0, 1)$ , the points  $(\beta_1, \beta_2)$  are aligned just above the bisector (gray dashed line). In other words, the values of  $\beta_2$  are highly correlated with the related values of  $\beta_1$ . We can derive an analytical approximation of the higher-order Heaps' exponents for this model, as we show in Sec. S3.2 of the SI. We obtain that the number of unique pairs for the UMT approximately grows as

$$D_2(t) \approx at^{\beta_2}, \quad \text{with} \quad \beta_2 = \beta_1 + \frac{c}{d + \log(t)}, \quad (1)$$

where  $a, c, d > 0$  depend on the parameters  $\rho$  and  $v$ , and  $\beta_1 = v/\rho$ . Although the predicted 2<sup>nd</sup>-order exponent is slightly higher than the 1<sup>st</sup>-order one, their difference just depends on the sequence length, and vanishes at larger times. Therefore, the difference between  $\beta_1$  and  $\beta_2$  observed in the simulations is only due to finite time effects, revealing how the UMT cannot reproduce the empirical patterns of Fig. 1. Due to the same reason, notice how the fitted values of  $\beta_1$  in the simulations of the UMT are lower than the asymptotically expected value of  $\beta_1 = v/\rho$  for high values of  $v$ , as also shown in Fig S6(a) in SI.

We repeat the analysis for two other generative models



**Figure 2. Higher-order Heaps' exponents in existing models.** Scatter plots of the (1<sup>st</sup>-order) Heaps' exponent  $\beta_1$  against the 2<sup>nd</sup>-order exponent  $\beta_2$  in: (a) the urn model with triggering (UMT), no semantic correlations ( $\eta = 1$ ), and  $\rho = 20$ ,  $\nu = 1, 2, \dots, 20$ ; (b) the urn model with semantic triggering (UMST) with  $\eta = 0.1$  and  $\rho = 4$ ,  $\nu = 1, 2, \dots, 20$ ; (c) the edge-reinforced random walk (ERRW) on a small-world network (average degree  $\langle k \rangle = 4$  and rewiring probability  $p = 0.1$ <sup>68</sup>) with edge reinforcement  $\rho$  ranging geometrically from 0.1 to 10. Each point refers to a different simulation of the related model, with colors representing the value of the free parameter (see color bar). Each panel also reports histograms of exponent distributions on the respective axes, and the bisector  $y = x$  (dashed gray line). All simulations have run for  $10^5$  time steps.

for discovery and exploration processes, i.e., the Urn Model with Semantic Triggering (UMST)<sup>1</sup> and the Edge-Reinforced Random Walk (ERRW)<sup>3</sup>, which have been proved to generate sequences obeying to the Heaps' law. These models share the same foundations of the UMT, but with some crucial differences. The UMST builds on top of the UMT introducing also semantic groups for colors. This addition effectively diminishes the probability to draw colors outside the semantic group of the last extracted color by a factor  $\eta \leq 1$ . The ERRW, instead, is formulated as a network exploration rather than a process of extractions from an urn. Instead of a sequence of extracted balls, the ERRW generates a sequence made of the nodes sequentially visited by the edge-reinforcing random walker over a weighted network, where the weight of the visited edges are reinforced by  $\rho$  when crossed. A full description of the models can be found in *Materials and Methods*.

We simulate the UMST with parameters  $\eta = 0.1$  (semantic parameter),  $\rho = 4$  (reinforcement parameter),  $\nu = 1, 2, \dots, 20$  (triggering parameter), while the simulations of the ERRW run over small-world networks<sup>69,70</sup> (with average degree  $\langle k \rangle = 4$  and rewiring probability  $p = 0.1$  following the procedure of Ref.<sup>68</sup>), with edge reinforcement  $\rho$  ranging from 0.1 to 10. Similarly to the UMT, we perform 100 simulations for each set of parameters, and report the results in Fig. 2(b-c). For both UMST and ERRW, we find that the values of  $\beta_2$  do not differ much from their corresponding value of  $\beta_1$ —as shown by the great proximity of the points  $(\beta_1, \beta_2)$  to the bisector. This means that also these models fail to reproduce the empirical variability of higher-order Heaps' exponents with respect to the 1<sup>st</sup>-order one. Moreover, we notice in (b) that for the UMST we only obtain exponents with either very low (up to 0.4) or very high (close to 1) values. We indeed see an abrupt

transition between these two extremes, with the model not able to cover the values in between, which are instead present in the empirical data reported in Fig. 1 (see also the relation with analytical results in Fig S6 in SI). Further simulations of both the UMST and ERRW with other sets of parameters are reported in Fig. S7 of the SI. Also in these other cases, these models are unable to generate sequences with  $\beta_2$  different from  $\beta_1$ .

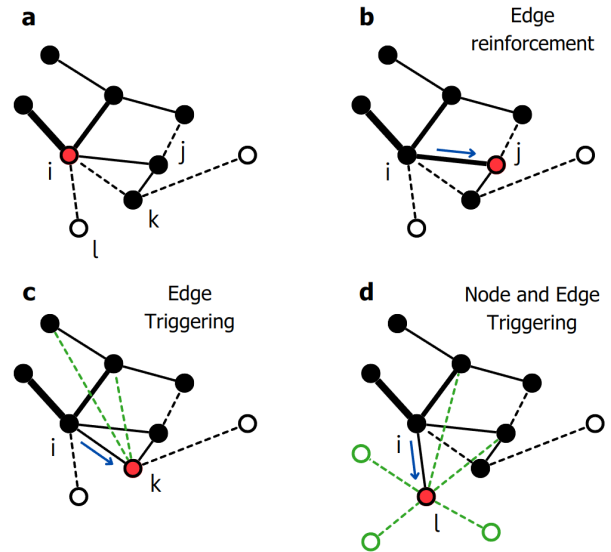
Overall, the analyses above indicate that, while the existing models of discovery and innovation dynamics are able to reproduce the empirically observed pace of discovery of new items (singletons) as captured by the 1<sup>st</sup>-order Heaps' law, they fail to capture the distributions of Heaps' exponents of higher order.

### The ERRWT: a model for higher-order Heaps' laws

In order to fill the gap between empirical observations and models, we introduce here a new model that can generate synthetic sequences with tunable discovery paces both at the first order and at higher orders. As for the previously discussed ERRW, our model is formulated in terms of network exploration. Namely, in the model: (i) the items to be explored correspond to the nodes of the network, (ii) the links between nodes represent semantic associations between items that one can use to move from one to another, and (iii) the exploration process is modelled as a random walk over the network and the sequence is obtained from the ordered list of visited nodes. Under these assumptions, the first visit of a node corresponds to a 1<sup>st</sup>-order novelty, while the first visit of a link corresponds to a 2<sup>nd</sup>-order novelty. Such definition can be straightforwardly extended to higher orders. In this manuscript, for simplicity, we limit our attention to the first two orders. The ERRW proposed in Ref.<sup>3</sup> consists of

a random walk on a network whose topology is fixed, i.e., the links cannot change in time, but the link weights can be modified by the passage of the random walker. By contrast, in our model not only the weights, but the entire network structure co-evolves with the exploration process, and new nodes and new links can be triggered. Thus, in analogy with the UMT<sup>1</sup>, we name the model *Edge-Reinforced Random Walk with Triggering* (ERRWT). More specifically, on top of the edge reinforcement mechanism of the ERRW, the model is based on two different triggering mechanisms that add new edges and new nodes every time a novelty appears. Similarly to the previous models analysed before, the first visit of a node triggers the expansion of the adjacent possible, as new nodes, neighbors of the discovered node, become now accessible. As an example of this mechanism, think for instance to the invention of the transistor, which made it possible to create mobile phones, among other things. Moreover, differently from the previous models, here the first visit of an edge is also considered a novelty, and as such, it triggers new edges. The idea is that whenever two elements are associated for the first time, new possible combinations involving one of these elements are then triggered. For instance, once photo-cameras and mobile phones were firstly combined, this association made clear that many more functions could be added to the latter, e.g., a music player, a game console, a GPS, etc. More formally, the model considers that the adjacent possible can be expanded at various orders, i.e., not just by introducing new nodes, but also by triggering new links.

The basic mechanisms of the ERRWT model are illustrated in Fig. 3. Suppose that at a given time  $t$ , the walker is located at node  $i$  of a network. At this point, some nodes and links, represented by full circles and solid lines in Fig. 3(a), have already been visited, while others, shown as empty circles and dashed lines, are part of the adjacent possible. In Fig. 3(b), the walker moves from node  $i$  to node  $j$ , crossing in this way an already explored link. Consequently, the weight of such link is increased by a positive quantity  $\rho$ , meaning that the association between the two nodes  $i$  and  $j$  becomes stronger and thus more likely to be used again. This is the same edge reinforcement mechanism adopted in the ERRW model<sup>3</sup>. In addition to this, if instead the walker moves from node  $i$  to node  $k$ , traversing an edge for the first time, as displayed in Fig. 3(c), this event is considered a 2<sup>nd</sup>-order novelty and triggers the creation of new edges. In particular,  $v_2 + 1$  new edges connecting node  $k$  to other already-visited nodes are created (green dashed lines). Finally, the third mechanism of the ERRWT model is analogous to the triggering mechanism of the UMT model. As illustrated in Fig. 3(d), when the walker moves from  $i$  to a node  $l$ , visiting node  $l$  for the first time, this event triggers the expansion of node  $l$ 's adjacent possible with the addition of new nodes and new links. Namely,  $v_1 + 1$  new nodes are added to the network and connected to the node  $l$  itself. In addition to this,  $v_2 + 1$  new links to already known elements are created, since whenever a node is explored for the first time, also the link leading to it is explored for the first



**Figure 3. The Edge-Reinforced Random Walk with Triggering (ERRWT) model.** An exploration process is modelled as a random walk on a growing weighted network. (a) At time  $t$ , the walker is at the red node  $i$ . Nodes that have been already visited by the walker are colored in black, in white those left to be visited. Similarly, traversed (old) and not-traversed (new) links are respectively depicted with continuous and dashed lines, whose widths represent their weights. At time  $t + 1$ , the walker can move to each of the neighbours of  $i$ , e.g. nodes  $j$ ,  $k$ , or  $l$ , with a probability proportional to the weight of the respective link. (b) If the walker moves to  $j$ , the weight of the link  $(i, j)$  is reinforced by  $\rho$  (Edge Reinforcement mechanism), but no new nodes or links are added to the network, since the link  $(i, j)$  is old; (c) if the walker moves to node  $k$ , since link  $(i, k)$  is new but node  $k$  is old, in addition to the edge reinforcement,  $v_2 + 1 = 2$  new edges (in green) between  $k$  and old nodes are added to the network (Edge Triggering mechanism); (d) finally, if the walker moves to  $l$ , since both the link  $(i, l)$  and the node  $l$  are new, in addition to the edge reinforcement and the edge triggering,  $v_1 + 1 = 3$  new nodes (in green) are added to the network and connected to  $l$  (Node and Edge Triggering mechanism).

time. More details about the ERRWT model can be found in *Materials and Methods*.

Balancing edge reinforcement and node and edge triggering mechanisms through the parameters  $\rho$ ,  $v_1$  and  $v_2$  of the ERRWT model, it is possible to control the pace of discovery of new nodes and edges, and consequently tuning the exponents of the 1<sup>st</sup>-order and the 2<sup>nd</sup>-order Heaps' law associated to the sequences produced by the model. To systematically explore this, we simulate the ERRWT model with parameters  $\rho = 10$ ,  $v_1 = 0, 1, \dots, 20$ , and  $v_2 = 0, 1, \dots, 2v_1$ , running 100 simulations for each set of parameters. Higher values of  $v_2$  have not been considered since they produce the same exponents as those for  $v_2 = 2v_1$ . In Fig. 4(a) we report the increase in the number of 1<sup>st</sup>-order and 2<sup>nd</sup>-order novelties (continuous lines) for a specific set of parameters as an example. The power-law fits (dashed lines) highlight that the Heaps' law is verified at the 2<sup>nd</sup> order too, leading to an increase of the exponent values (from  $\beta_1 = 0.56$  to  $\beta_2 = 0.87$ ). The relation between the different orders is explored in the scatter plot between the 1<sup>st</sup>- and 2<sup>nd</sup>-order Heaps' exponents reported in Fig. 4(b). Each point refers to a different simulation, and we use the color to indicate the value of the parameter  $v_1$  used (see color bar). We notice that the ERRWT model can produce a wide range of values for the exponents at both orders, and that the 2<sup>nd</sup>-order exponents are not trivially correlated to the 1<sup>st</sup>-order ones, as it happened in the models considered in the previous section. This is even more clear when we look at Fig. 4(c), where the Heaps' exponents are averaged across simulations for each set of parameters. Each curve in the figure refers to a different value of  $v_1$ , with  $v_1$  increasing from 1 to 20 from bottom left to top right of the panel. The color represents instead different values of the parameter  $v_2$  from 0 to  $2v_1$ . For reference, we also flag using a red dot the pair of exponents related to the parameters used in Fig. 4(a). We can immediately notice how the 1<sup>st</sup>- and 2<sup>nd</sup>-order Heaps' exponents increase as  $v_1$  gets larger. More interestingly, we observe the combined role of the two parameters. For each curve, by increasing  $v_2$ , therefore triggering new links in the network, the difference between  $\beta_1$  and  $\beta_2$  becomes larger, and the point  $(\beta_1, \beta_2)$  moves away from the bisector, in a way that depends on the specific value of  $v_1$ . In particular, for low values of  $v_1$ , the curves are almost vertical, with only  $\beta_2$  increasing. Instead, for higher values of  $v_1$ , especially when  $v_1 \geq \rho$ , an increase of  $v_2$  produces a decrease of  $\beta_1$ , while the value of  $\beta_2$ , which is close to its upper bound value 1, does not change. Intuitively, this happens because the creation of more and more new links between explored nodes increases the chance to exploit nodes already discovered, while still exploring never traversed links.

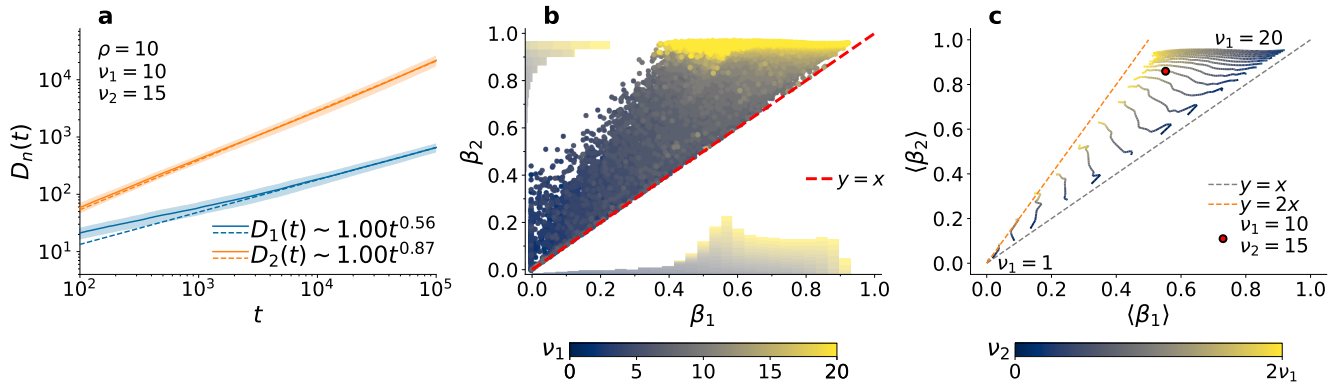
It is also possible to perform an analytical investigation of a simplified version of the ERRWT model, which leads to results in agreement with the simulations (see Sec. S4 in SI). In particular, for such a simplified model, we can prove that the values of the asymptotic Heaps' exponents  $\beta_1$  and  $\beta_2$  depend on the two ratios  $v_1/\rho$  and  $v_2/\rho$ . Moreover, we

find that, for  $v_1/\rho > 1$ , the 2<sup>nd</sup>-order Heaps' exponent is asymptotically equal to 1, while the 1<sup>st</sup>-order one depends on  $v_1/v_2$ , as seen in Fig. 4(c). Finally, the exponents are asymptotically bounded by  $\beta_1 \leq \beta_2 \leq 2\beta_1$ , as also observed in the simulations in Fig. 4(c). This also explains why the exponents do not change when we increase  $v_2$  above  $2v_1$ .

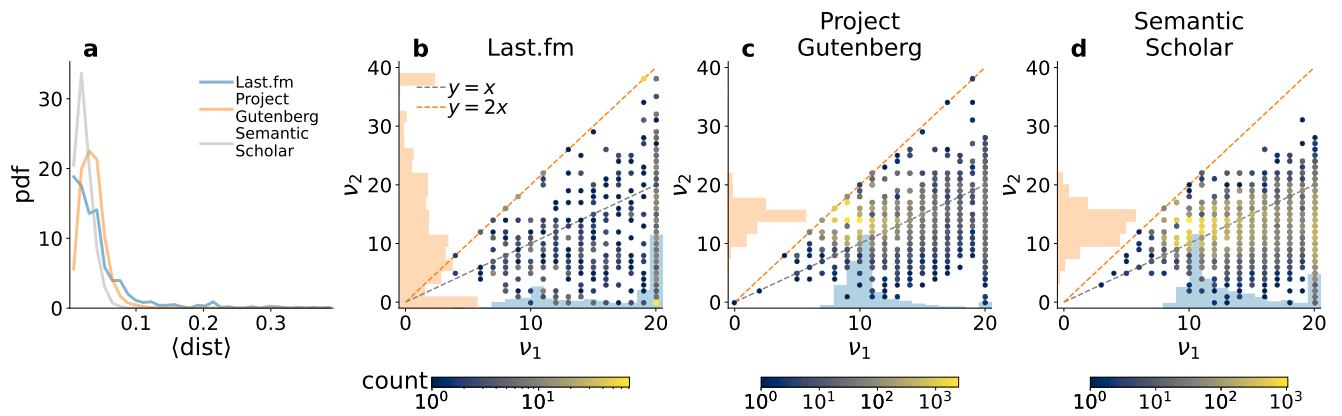
### Comparison between ERRWT and real-world data

To show that the ERRWT model is able to reproduce the properties observed in real-world processes, we now fit the parameters of the model to the three data sets analyzed (Last.fm, Project Gutenberg and Semantic Scholar). Given an empirical sequence and its pair of 1<sup>st</sup>- and 2<sup>nd</sup>-order Heaps exponents  $(\beta_1, \beta_2)$ , we compute the Euclidean distance between the pair  $(\beta_1, \beta_2)$  and each of the pairs of exponents  $(\beta'_1, \beta'_2)$  obtained by simulating the ERRWT model using the sets of parameters considered in the previous section. We then select the best model parameters by minimizing the average distance over 100 simulations for each set, and repeat the procedure for all the sequences of the three data sets. Figure 5(a) shows the probability density distribution of the distances between the empirical sequences and the simulations of the best-performing ERRWT model. Notice how these distances are almost all below 0.1, that is the uncertainty we expect on the values of the parameters. Indeed, being  $v_1, v_2$  integers and  $\rho = 10$ , the maximum precision we can gain on the estimate of the best parameters is about  $1/\rho = 0.1$ . The percentage of sequences with higher distance than this threshold is 7.67%, 0.73%, and 0.05% for Last.fm, Project Gutenberg, and Semantic Scholar, respectively. The scatter plots of the best-fitted parameters  $v_1$  and  $v_2$  for the three data sets are shown in Fig. 5(b-d). The colors here indicate the number of empirical sequences which are best represented by each pair of parameters  $v_1$  and  $v_2$ . We notice that most of the sequences of Last.fm are characterized by relatively large values of  $v_1$ . Since  $v_1$  is related to the triggering of new nodes, this result indicates that the discovery of a new song exposes the user to a large variety of related songs, previously not accessible, which can now be discovered. Conversely, the parameter  $v_2$ , which controls the triggering of new edges between already existing items in the model, takes values in a larger range, predominantly skewed towards the lower end. This suggests that, once a new association between two songs is established by a user, there is a high probability that the same association will be repeated over and over. Consequently, the user will preferably listen to songs in a similar order, instead of creating new associations. We point out that we cannot distinguish if this is due to individual preferences or is pushed by the presence of recommender systems in music listening platforms. In the case of Project Gutenberg, most sequences have  $v_2 > v_1$ . This implies that writers tend to frequently generate new word associations instead of using words never used before in the text, highlighting the incredible variety of expressions we can make by combining a limited set of words. Finally, Semantic Scholar exhibits values of  $v_1$  and  $v_2$  similar to those found





**Figure 4. Higher-order Heaps' exponents in the ERRWT model.** (a) Average number  $D_n(t)$  of novelties of order  $n$ , with  $n = 1$  and  $2$ , as a function of the sequence length  $t$  for simulations of the ERRWT model with parameters  $\rho = 10$ ,  $v_1 = 10$ ,  $v_2 = 15$ , and fit of the associated Heaps' laws (dashed lines), with estimated exponents shown in the legend. Shaded areas represent one standard deviation above and below the average. (b) Scatter plot between the (standard) Heaps' exponent  $\beta_1$  and the 2<sup>nd</sup>-order exponent  $\beta_2$ . Each point refers to a different simulation of the model, with colors representing the corresponding value of the parameter  $v_1$  ranging from 0 to 20 (see color bar), while  $\rho = 10$  and  $v_2 = 0, \dots, 2v_1$ . (c) Variation of the average  $n$ <sup>th</sup>-order Heaps' exponents  $\beta_n$ , with  $n = 1, 2$ . Each curve refers to a different value of  $v_1$ , increasing from 1 to 20 from bottom left to top right, while the color represents the value of  $v_2$  (see color bar). The set of parameters used in (a) is here highlighted in with a red dot.



**Figure 5. Fitting the ERRWT model to real-world data sets.** (a) Distribution of the average distance between the pair of exponents  $(\beta_1, \beta_2)$  of a real sequence and the pair  $(\beta'_1, \beta'_2)$  obtained by the best fitting ERRWT model. (b-c) Scatter plots of the best-fitted parameters  $v_1$  and  $v_2$  of the model across the sequences of the three data sets, respectively Last.fm (b), Project Gutenberg (c), and Semantic Scholar (d). The color of a point refers to the number of sequences with that pair of parameters in the best fitting ERRWT model (see color bar).

in the Project Gutenberg data set. However, some sequences of Semantic Scholar have a relatively high value of  $v_1$  with respect to  $v_2$ . This is an indication that, when choosing words for titles, authors tend to use more original words, while the pace of creation of new word associations remains similar.

## Discussion

The extraction of the Heaps' exponent from empirical sequences has recently allowed to characterize the pace at which discoveries occur in different contexts<sup>1-3,71,72</sup>. However, there is more and more evidence that discoveries are often made from novel combinations of already known elements<sup>46,48,49,55,56</sup>. In this manuscript we have proposed to explore higher-order Heaps' laws and to extract higher-order Heaps' exponents as a way to characterize novel combinations in an exploration process. The key idea is to look at novelties not only as discoveries of single items, but also as the first appearances of new combinations of two or more items. More precisely, the Heaps' exponent of order  $n$  measures the rate of discovery of novelties of order  $n$ , i.e., combinations of  $n$  items. Notice that our approach differs and complements other measures of the pace of discovery recently proposed. For example, the authors of Refs.<sup>8,73</sup> have investigated the number of all possible valid combinations that can be created by using the elements acquired so far as a proxy of the level of innovation of a given system. In this way, the potential for new discoveries is accounted for, rather than the actual number of novel combinations observed in a system and their rate of appearance.

As we have shown through the analysis of empirical sequences of music listening records, higher-order Heaps' exponents can be used to further classify users of Last.fm with the same rate of discovery of new songs or new artists. The higher-order Heaps' exponent can indeed tell apart different ways to explore the same set of songs in terms of number of different associations of consecutive pairs or triples of songs. Analogously, we found that higher-order Heaps' exponents can uncover different patterns in the use of words in different texts. Titles of peer-reviewed papers published in scientific journals show more creative combinations of words, than the texts of narrative books. They indeed exhibit many more new  $n$ -grams, even if the total set of words used is similar in length. Overall, our analysis shows that the space of possibilities grows in a complex way, which does not depend solely on the balance between old items to exploit and new ones to explore, but also on the structure of their associations. Notice, however that in our framework we have considered all associations of consecutive items in a sequence as possible discoveries. This is certainly a strong assumption, which might not always be valid. For instance, in the context of a written text, the last word of a paragraph and the first word of the following paragraph, are not necessarily related to each other and could be discarded from the sequence of pairs to be analyzed. The problem can be solved by filtering out all such cases, as we have done for books in Gutenberg in Fig. S5, obtaining similar

results in terms of 1<sup>st</sup>-order and higher-order Heaps' exponents. In other contexts, it might be necessary to better tailor the precise definition of "novel combination" according to the nature of the sequence being analysed and the underlying research question.

We have then focused our attention in understanding the underlying mechanisms that can trigger higher-order novelties. We have proposed a new modelling framework, the ERRWT which takes into account not only the exploration rate of new items, but also the propensity to explore the same content in a more creative way. Considering that pairs of items can be seen as the links of a network, the model is based on a process of network exploration and on the co-evolution of the network structure with the dynamics of the exploration process. The model considers a reinforcement of the visited links and the triggering of new nodes and links whenever new nodes or links are explored. Not only the ERRTW model is able to reproduce the higher-order Heaps' exponents extracted from real data, but also provides a new intuition of how the space of possibilities grows over time, shedding light on the underlying mechanism in which novel elements and combinations emerge.

We acknowledge there are various ways in which our model can be improved and generalized. For example, future work should investigate the interplay between initial knowledge, either of the individual or of a group, and the pace of discovery at various orders during the exploration process, or the influence of recommendation algorithms. In our model, we have supposed that the links all start with the same weight, which can be a too strong assumption in certain contexts. Moreover, we have assumed to trigger new links with a uniform probability. It would be interesting to study cases in which the space has some preferential pathways, for example represented by an underlying network structure. This could be implemented in our model by limiting the addition of new links to only those permitted by an underlying network given as an input to the model. Alternatively, more complex ways to trigger edges, such as preferential attachment mechanisms, could be considered<sup>74,75</sup>. Finally, we have not considered the presence of semantic correlations in the temporal sequence of visited items, which can be a consequence of the interplay between the network topology and a predisposition to move within items semantically close to the recent ones, reinforcing a clustered structure. It would indeed be interesting to use higher-order Heaps' exponents and the ERRWT model to study phenomena related to waves of novelties<sup>2</sup> and popularity<sup>76</sup>. Moreover, the ERRWT model could be extended to a multi-agent model to study how different agents would cooperate and diffuse knowledge<sup>14,37</sup>, also taking into account the presence of a limited attention capacity and memory that could influence the rise and fall of popular items<sup>77</sup>. We believe that our model can be directly used to answer these questions and, more in general, to better understand the fundamental mechanism behind innovation and creativity.

## Methods

### Data

In this work we consider three different data sets: music listening records (*Last.fm*), books (*Project Gutenberg*), and scientific articles (*Semantic Scholar*).

*Last.fm* is a digital platform for music born in 2002, famous for logging all listening activities of its users, providing both personal recommendations and a space to interact with other users interested in music<sup>78</sup>. In this manuscript, we use a data set presented in Ref.<sup>79</sup> and available at Ref.<sup>80</sup>. This data set has been obtained and used according to terms and conditions of *Last.fm*<sup>78</sup>, and can only be used for non-commercial use<sup>80</sup>. It contains all listening records of about 1000 users. In order to have sequences long enough for statistically relevant fits, only users with more than 1000 logs have been retained. The final data set contains 890 users having a median number of listened records of 13 985. Each record contains the timestamp at which a user listened to a given song. In the database, each song is associated to a title, the artist's name and a unique MusicBrainz Identifier (MBID), which can be used to obtain additional metadata<sup>81</sup>. Using this information, we are able to create, for each user, a temporally ordered sequence of songs together with the associated sequence of artists. It is worth noting that the behavior of each user might be influenced by additional factors, such as recommendation algorithms. While the specifics of these procedures are not known in details, we expect that they would not drastically alter our main findings<sup>2</sup>. They may, for example, impact only the numerical values obtained, without affecting the fundamental mechanisms captured by our modelling approach.

*Project Gutenberg* is an open access text corpus containing more than 50 000 books of different nature<sup>82</sup>. This corpus is made of public domain books, with expired copyrights, which can therefore be disseminated freely and legally. Here, we make use of the Standardized Project Gutenberg Corpus<sup>83</sup>, which allows to download and process an updated version of the open corpus. Using Google's Compact Language Detector 3 (cld3 package in Python), we filter out all non-English texts. We then discard all texts with less than 1000 words, retaining a total of 19 637 books with a median number of 50 726 words. A sequence of events for each book is hence created with the lemmatized words, disregarding punctuation and putting all characters in lower case. We also extract stems from each word using the English Snowball stemmer<sup>84</sup>—a more accurate extension of the Porter stemmer<sup>85</sup>—, which is not as aggressive as the Lancaster stemmer<sup>86</sup>.

*Semantic Scholar* is a recent project with the scope of facilitating scientific analysis of academic publications<sup>87</sup>. It provides monthly snapshots of research papers published in all fields, publicly and freely accessible through the *Semantic Scholar Academic Graph* (S2AG, pronounced “stag”)<sup>88</sup>. This database (1<sup>st</sup> Jan. 2022 snapshot) contains about 203.6M papers, 76.4M authors, and 2B citations, obtained in accordance with the project terms and conditions<sup>89</sup>. It also classifies each

paper into one or more fields of study<sup>90</sup>, for a total of 19 different fields. For simplicity, we associate each paper to its first (and most relevant) field of study. To create the sequences to analyze, for each field we consider the first 1000 journals in terms of number of English papers. Then, for each journal, we order the published papers based on the respective year of publication, volume, issue, and first page. When some of this information is not available, the Semantic Scholar unique ID of the paper is also used in the ordering process. Thus, for each paper, we extract and lemmatize their title, similarly to what done for the Project Gutenberg. Finally, a sequence of events is created for each selected journal, concatenating the lemmatized words in the titles of each paper in their temporal order, for a total of 19 000 sequences with median length of 9 114.5. Associated to each sequence, we also consider the sequence of stemmed words for further analysis, similarly to the Project Gutenberg corpus..

### Power-law fit

Fundamental for the estimation of the higher-order Heaps' exponent of a sequence is the power-law fitting procedure for the number of novel  $n$ -tuples  $D_n(t)$  as a function of the sequence length  $t$ , with  $n \geq 1$ . The sequences analyzed in this manuscript come from very different contexts, from empirical data sets to model simulations. We thus need to take into consideration all those cases that show a transient regime—whose length might also depend on the system structure<sup>37</sup>—in which the pace of discovery can fluctuate before reaching its stationary value. Therefore, we fit each sequence according to the following procedure. To reduce computational times, we first logarithmically sample 1000 real points in the range  $[1, T]$ , where  $T$  is the length of the sequence. Considering their integer part, we discard all the duplicates that may be produced when some sampled points differ only for the decimal part. We thus obtain a set of  $k$  integer times  $\{t_i\}_{i=1,\dots,k}$  between 1 and  $T$ . Due to the removal of duplicates,  $k$  can be equal or smaller than 1000. If  $T \geq 1000$ , that is the case of all sequences analyzed in this manuscript, then this process results in  $k \geq 424$  points. Taking into account that the associated sequence of  $n$ -tuples has length  $T - n + 1$ , we thus consider the points  $\{(t_i - n + 1, D_n(t))\}_{i=1,\dots,k}$  in logarithmic scale, i.e.,

$$(x_i, y_i) = (\log_{10}(t_i - n + 1), \log_{10}(D_n(t))), \quad (2)$$

with  $i = 1, \dots, k$ . In order to neglect the initial transient regime, but still have enough points for a sufficiently significant fit, we select only the last 100 of such points. We hence look for the best fit of  $\{(x_i, y_i)\}_{i=k-100+1,\dots,k}$  by optimizing the linear function  $y = a + bx$ , with  $a \geq 0$ , using the tool `curve_fit` of the Python package `Scipy`<sup>91</sup>. The constraint  $a \geq 0$  is necessary to avoid that the fitted Heaps' exponent is greater than 1, which could happen when the initial transient regime differs significantly from the asymptotic one and could hence produce wrong fits. Finally, if  $\bar{a}$  and  $\bar{b}$  are the best parameters, then the power-law fit of the Heaps' law is  $D_n(t) \approx 10^{\bar{a}} t^{\bar{b}}$ ,

that is, the  $n^{\text{th}}$ -order Heaps' exponent is approximated by the slope  $\bar{b}$  of the fit.

### Urn Model with (Semantic) Triggering

The Urn Model with Triggering (UMT) is a random generative model for discovery processes, producing a sequence of extractions of balls of various colors from an urn. First introduced in Ref.<sup>1</sup>, it successfully reproduces the main features of empirical discovery processes<sup>1,6,59,92</sup>. The UMT can be thought as an extension of Pólya Urn processes<sup>33,34,93–95</sup> that includes the concept of *adjacent possible*<sup>26</sup> in the way a novelty can trigger further ones<sup>2,32</sup>. Differently from the classic urn of Pólya in which only balls of existing colors can be added to the urn, the UMT features a growing number of colors, that is, the set of possible events expands together with the exploration process. It is hence the process itself that shapes the content of the urn by reinforcing elements already discovered and adding new possibilities.

Supposing that the urn initially contains  $N_0$  balls of different colors, the UMT works as follows. At each discrete time-step  $t$ , a ball is randomly drawn from the urn with uniform probability, and its color is marked in a temporally-ordered sequence of events  $\mathcal{S}$  at position  $t$ . The extracted ball is then put back in the urn together with other  $\rho$  copies of the same color, in a *rich-get-richer* manner<sup>74</sup>. This mechanism ensures that frequently adopted items, visited places, or exploited concepts will be more and more likely to be adopted, visited, or exploited in the future. Furthermore, if the color of the extracted ball has never appeared before in  $\mathcal{S}$ , this event is considered to be a novelty. As a consequence it triggers new possibilities, represented by the addition of  $\nu + 1$  balls—each of a new different color—into the urn. This triggering mechanism thus ensures the expansion of the space of possibilities.

In a different version of the model, the Urn Model with Semantic Triggering (USMT), the sequences produced contain semantic correlations between consecutive extractions, as seen in the data<sup>1</sup>. The UMST works similarly to the UMT, but with the introduction of semantic groups for colors. In particular, at each triggering event, supposing that the triggering color belongs to the group  $A$ , the new  $\nu + 1$  colors are assigned to a common new group  $B$ , semantically related to the triggering color. Therefore, a color  $i$  of label  $A$  is semantically related to all other colors of label  $A$  (siblings), the color that triggered the addition of  $A$  in the urn (parent), as well as all colors of label  $B$  that have been triggered by  $i$  (children). Taking this into consideration, at each extraction, the probability to extract each color changes depending on a fixed parameter  $\eta \in [0, 1]$ . A ball has weight 1 if its color is semantically related to the one extracted on the previous time-step, otherwise it has weight  $\eta$ . Notice that we can recover the original UMT by simply considering  $\eta = 1$ .

Finally, as shown in Ref.<sup>1</sup>, the effect of  $N_0$  is negligible at large times. For simplicity, we thus consider  $N_0 = 1$  in our simulations of both UMT and UMST.

### Edge-Reinforced Random Walk

Given a weighted connected graph  $G = (\mathcal{V}, \mathcal{E})$  with  $N = |\mathcal{V}|$  vertices (nodes) and  $M = |\mathcal{E}|$  edges (links), the Edge-Reinforced Random Walk (ERRW) is a dynamical process that reinforces the weights of the visited edges in  $\mathcal{E}$ , leading to Heaps' laws<sup>3</sup>. The weights of the links in the network quantify the strength of the relationship among nodes, and are encoded in a time-varying adjacency matrix  $W^t \equiv \{w_{ij}^t\}$ . This matrix features non-zero entries  $w_{ij}^t$  when at time  $t$  the link connecting node  $i$  and node  $j$  is different from zero. Let us assume that at time  $t = 0$  each link  $(i, j) \in \mathcal{E}$  has weight  $w_{ij}^0 = 1$ , while all other weights are set to zero. At each time step, a walker at node  $i$  walks to a neighboring node  $j$  with a probability that is proportional to the weight of the outgoing links, i.e.,  $\mathbb{P}(i \rightarrow j) = w_{ij}^t / \sum_l w_{il}^t$ . After moving to the randomly chosen node  $j$ , a reinforcement  $\rho$  is added to the weight of the traversed edge  $(i, j)$ , i.e.,  $w_{ij}^{t+1} = w_{ij}^t + \rho$ . Starting from an underlying structure given by the graph  $G$ , the ERRW can generate sequences of visited nodes with a tunable pace of discovery obtained by properly calibrating the reinforcement parameter  $\rho$ <sup>3</sup>. Because of the interplay between structure and dynamics, different structures might require different values of the reinforcement parameter to reach the same pace of discovery. For example, higher values of  $\rho$  must be chosen for a denser graph. This is similar to what happens in the UMT, in which we need higher values of the reinforcement parameter  $\rho$  to obtain the same pace of discovery as we increase the triggering parameter  $\nu$ .

### Edge-Reinforced Random Walk with Triggering

In this manuscript we propose a generative model of a discovery process based on the exploration of a growing network, i.e., the Edge-Reinforced Random Walk with Triggering (ERRWT), which can be considered as a UMT-inspired extension of the ERRW model. For this model, any initial connected network  $G^0 = (\mathcal{V}^0, \mathcal{E}^0)$  with  $N^0 = |\mathcal{V}^0| \geq 1$  nodes and  $M^0 = |\mathcal{E}^0|$  links can be used. Let us suppose that the nodes of the graph are indexed, that is,  $\mathcal{V}^0 = \{1, 2, \dots, N_0\}$ . Similarly to the ERRW model, we assume that all initial links  $(i, j) \in \mathcal{E}^0$  have weight  $w_{ij}^0 = 1$ . The initial node to start the exploration process is randomly selected from  $\mathcal{V}^0$ . We let the graph evolve during the process, adding new nodes and links. Let  $G^t = (\mathcal{V}^t, \mathcal{E}^t)$  be the graph at time  $t$ . The structure of the growing network is encrypted in the time-varying weighted adjacency matrix  $W^t \equiv \{w_{ij}^t\}$ , where  $w_{ij}^t$  represents the weight of the link  $(i, j)$  at time  $t$ . We assume here that  $G^t$  is an undirected graph, so the matrix  $W^t$  is symmetric, and any variation of  $w_{ij}^t$  affects  $w_{ji}^t$  too. Supposing that at time  $t$  the ERRWT is positioned on node  $i$  of  $G^t$ , the model obeys to the following rules.

- *Choice of next node.* The ERRWT randomly moves to a neighbouring node  $j$  of the current node  $i$ . The probability to move to node  $j$  depends on the weight of

the outgoing links of  $i$ , i.e.,

$$\mathbb{P}(i \rightarrow j) = \frac{w_{ij}^t}{\sum_l w_{il}^t}. \quad (3)$$

- *Edge reinforcement.* The weight of the chosen edge  $(i, j)$  is reinforced by  $\rho$ , that is,

$$w_{ij}^{t+1} = w_{ij}^t + \rho. \quad (4)$$

- *Edge triggering.* If the walker never traversed the chosen edge  $(i, j)$  before, i.e., it is a new link, then  $v_2 + 1$  new possible links are added to the network. These links are connections of unitary weight between  $j$  and previously visited nodes  $l = l_1, \dots, l_{v_2}$  in  $\mathcal{V}^t$ , for which the link  $(j, l)$  has never been traversed by the walker. If one of these edges already exists in the space of possibilities, its weight is reinforced by one more unit, otherwise, it is added to  $\mathcal{E}^{t+1}$ . In other words, we have

$$w_{jl}^{t+1} = w_{jl}^t + 1, \quad l = l_1, \dots, l_{v_2} \mid l \text{ old}, (j, l) \text{ new}. \quad (5)$$

- *Node triggering.* If the walker never visited the chosen node  $j$  before, i.e., it is a new node, then  $v_1 + 1$  new possible nodes are added to the network; these are connected to node  $j$  with unitary weights. Mathematically, we have

$$\begin{aligned} \mathcal{V}^{t+1} &= \mathcal{V}^t + \{l\}_{l=|\mathcal{V}^t|+1, \dots, |\mathcal{V}^t|+v_1+2} \\ w_{jl}^{t+1} &= 1, \quad l = |\mathcal{V}^t| + 1, \dots, |\mathcal{V}^t| + v_1 + 2. \end{aligned} \quad (6)$$

Notice that if the chosen node  $j$  is new, then also the traversed edge  $(i, j)$  is necessarily new as well. Therefore, in this case there is also a triggering of  $v_2 + 1$  edges from  $j$  to other previously visited nodes, as described before.

Finally, in this manuscript, we let the initial graph  $G_0$  be a small graph that emulates the triggering mechanism introduced, shown in Fig. S11 in SI. This is a regular tree with branching parameter  $v_1 + 1$  and two levels, where only the leaves are considered new, since all other nodes have already triggered. In other words, a root node has triggered  $v_1 + 1$  nodes connected to it, and again these nodes have also triggered each  $v_1 + 1$  other nodes. Therefore, we initially suppose that the triggered nodes, which are  $v_1 + 2$  in number, are all known to the walker at the start of the simulation, and do not trigger again when later explored. Moreover, we assume that all links are new to the walker and have unitary weight. This initialization makes sure that in the initial stages of the simulation there are enough possible links between already known nodes. As we show in Sec. S4 in SI where we test different initial graphs, the initialization procedure only affects thermalization times, and becomes irrelevant asymptotically.

### Data availability

The data used in this manuscript is publicly available at Refs.<sup>80,83,88</sup>, and has been obtained and used according to their terms and conditions.

### Code availability

All the code used to download, process and analyse the data and the models can be found at Ref.<sup>96</sup>.

### References

1. Tria, F., Loreto, V., Servedio, V. D. P. & Strogatz, S. H. The dynamics of correlated novelties. *Scientific Reports* **4**, 1–8 (2014).
2. Monechi, B., Ruiz-Serrano, Á., Tria, F. & Loreto, V. Waves of novelties in the expansion into the adjacent possible. *PLoS one* **12**, e0179303 (2017).
3. Iacopini, I., Milojević, S. & Latora, V. Network dynamics of innovation processes. *Physical Review Letters* **120**, 048301 (2018).
4. North, M. *Novelty: A history of the new* (University of Chicago Press, 2013).
5. Rzhetsky, A., Foster, J. G., Foster, I. T. & Evans, J. A. Choosing experiments to accelerate collective discovery. *Proceedings of the National Academy of Sciences U.S.A.* **112**, 14569–14574 (2015).
6. Loreto, V., Servedio, V. D. P., Strogatz, S. H. & Tria, F. Dynamics on expanding spaces: modeling the emergence of novelties. In *Creativity and universality in language*, 59–83 (Springer, 2016).
7. Puglisi, A., Baronchelli, A. & Loreto, V. Cultural route to the emergence of linguistic categories. *Proceedings of the National Academy of Sciences U.S.A.* **105**, 7936–7940 (2008).
8. Fink, T., Reeves, M., Palma, R. & Farr, R. Serendipity and strategy in rapid innovation. *Nature communications* **8**, 1–9 (2017).
9. Saracco, F., Di Clemente, R., Gabrielli, A. & Pietronero, L. From innovation to diversification: a simple competitive model. *PLoS One* **10**, e0140420 (2015).
10. Abbas, K. *et al.* Popularity and novelty dynamics in evolving networks. *Scientific Reports* **8**, 1–10 (2018).
11. Rodi, G. C., Loreto, V. & Tria, F. Search strategies of wikipedia readers. *PLoS One* **12**, e0170746 (2017).
12. Sreenivasan, S. Quantitative analysis of the evolution of novelty in cinema through crowdsourced keywords. *Scientific Reports* **3**, 1–11 (2013).
13. Iacopini, I. & Latora, V. On the dual nature of adoption processes in complex networks. *Frontiers in Physics* **9**, 604102 (2021).
14. Di Bona, G. *et al.* Social interactions affect discovery processes. *arXiv preprint arXiv:2202.05099* (2022).
15. Yule, G. U. A mathematical theory of evolution, based on the conclusions of dr j. c. willis, f.r.s. *Philosophical Transactions of the Royal Society of London* **213**, 21–87 (1924).
16. Simon, H. A. On a class of skew distribution functions. *Biometrika* **42**, 425–440 (1955).
17. Herdan, G. *Type-token Mathematics: A Textbook of Mathematical Linguistics*, vol. 4 (Mouton en company, 1960).
18. Heaps, H. S. *Information Retrieval: Computational and Theoretical Aspects* (Academic Press, Inc., Orlando, FL, USA, 1978).
19. Lü, L., Zhang, Z.-K. & Zhou, T. Zipf’s law leads to heaps’ law: Analyzing their relation in finite-size systems. *PLoS one* **5**, e14139 (2010).
20. Estoup, J.-B. *Gammes sténographiques* (Institut sténographique de France, 1916).
21. Zipf, G. K. Relative frequency as a determinant of phonetic change. *Harvard studies in classical philology* **40**, 1–95 (1929).
22. Zipf, G. K. The psychobiology of language (1935).
23. Zipf, G. K. Human behavior and the principle of least effort. cambridge, (mass.): Addison-wesley, 1949, pp. 573. *Journal of Clinical Psychology* **6**, 306–306 (1950).
24. De Marzo, G., Gabrielli, A., Zaccaria, A. & Pietronero, L. Dynamical approach to zipf’s law. *Physical Review Research* **3**, 013084 (2021).
25. Taylor, L. R. Aggregation, variance and the mean. *Nature* **189**, 732–735 (1961).

26. Kauffman, S. A. Investigations: The nature of autonomous agents and the worlds they mutually create. In *SFI working papers* (Santa Fe Institute, 1996).
27. Packard, N. H. Adaptation toward the edge of chaos. *Dyn. Patterns Complex Syst.* **212**, 293 (1988).
28. Langton, C. Computation at the edge of chaos: Phase transition and emergent computation (1990).
29. Langton, C., Taylor, C., Farmer, J. & Rasmussen, S. *Artificial Life II* (Avalon Publishing, 2003).
30. Bak, P. & Sneppen, K. Punctuated equilibrium and criticality in a simple model of evolution. *Physical Review Letters* **71**, 4083 (1993).
31. Armano, G. & Javarone, M. A. The beneficial role of mobility for the emergence of innovation. *Scientific Reports* **7**, 1–8 (2017).
32. Gravino, P., Monechi, B., Servedio, V. D. P., Tria, F. & Loreto, V. Crossing the horizon: exploring the adjacent possible in a cultural system. In *Proceedings of the Seventh International Conference on Computational Creativity, Paris*. (2016).
33. Eggenberger, F. & Pólya, G. über die statistik verketteter vorgänge. *ZAMM - Journal of Applied Mathematics and Mechanics / Zeitschrift für Angewandte Mathematik und Mechanik* **3**, 279–289 (1923).
34. Hoppe, F. M. Pólya-like urns and the ewens' sampling formula. *Journal of Mathematical Biology* **20**, 91–94 (1984).
35. Ubaldi, E., Burioni, R., Loreto, V. & Tria, F. Emergence and evolution of social networks through exploration of the adjacent possible space. *Communications Physics* **4**, 1–12 (2021).
36. Marzo, G. D., Pandolfelli, F. & Servedio, V. D. P. Modeling innovation in the cryptocurrency ecosystem. *Scientific Reports* **12**, 12942 (2022).
37. Iacopini, I., Di Bona, G., Ubaldi, E., Loreto, V. & Latora, V. Interacting discovery processes on complex networks. *Physical Review Letters* **125**, 248301 (2020).
38. Lakoff, G. & Johnson, M. *Metaphors we live by* (University of Chicago press, 2008).
39. Pinker, S. *The language instinct: How the mind creates language* (Penguin uK, 2003).
40. Boyd, D. M. & Ellison, N. B. Social network sites: Definition, history, and scholarship. *Journal of computer-mediated Communication* **13**, 210–230 (2007).
41. Tufekci, Z. Big questions for social media big data: Representativeness, validity and other methodological pitfalls. In *Proceedings of the international AAAI conference on web and social media*, vol. 8, 505–514 (2014).
42. Cope, D. *The Algorithmic Composer* (2000).
43. Market, E. & Papavasiliou, F. N. V (d) j recombination and the evolution of the adaptive immune system. *PLoS Biology* **1**, e16 (2003).
44. Jones, J. M. & Gellert, M. The taming of a transposon: V (d) j recombination and the immune system. *Immunological Reviews* **200**, 233–248 (2004).
45. Fortunato, S. *et al.* Science of science. *Science* **359** (2018).
46. Uzzi, B., Mukherjee, S., Stringer, M. & Jones, B. Atypical combinations and scientific impact. *Science* **342**, 468–472 (2013).
47. Ke, Q., Ferrara, E., Radicchi, F. & Flammini, A. Defining and identifying sleeping beauties in science. *Proceedings of the National Academy of Sciences U.S.A.* **112**, 7426–7431 (2015).
48. Wang, J., Veugelers, R. & Stephan, P. Bias against novelty in science: A cautionary tale for users of bibliometric indicators. *Research Policy* **46**, 1416–1436 (2017).
49. Fontana, M., Iori, M., Montobbio, F. & Sinatra, R. New and atypical combinations: An assessment of novelty and interdisciplinarity. *Research Policy* **49**, 104063 (2020).
50. Wu, L., Wang, D. & Evans, J. A. Large teams develop and small teams disrupt science and technology. *Nature* **566**, 378–382 (2019).
51. Alvarez-Rodriguez, U. *et al.* Evolutionary dynamics of higher-order interactions in social networks. *Nature Human Behaviour* **5**, 586–595 (2021).
52. Schumpeter, J. A. *et al.* *Business cycles*, vol. 1 (McGraw-Hill New York, 1939).
53. Schumpeter, J. A. *Capitalism, socialism and democracy* (routledge, 2013).
54. McNerney, J., Farmer, J. D., Redner, S. & Trancik, J. E. Role of design complexity in technology improvement. *Proceedings of the National Academy of Sciences U.S.A.* **108**, 9008–9013 (2011).
55. Abbasiharofteh, M., Kogler, D. F., Lengyel, B. *et al.* Atypical combination of technologies in regional co-inventor networks. *Papers in Evolutionary Economic Geography (PEEG)* **20** (2020).
56. Lambert, B. *et al.* The pace of modern culture. *Nature Human Behaviour* **4**, 352–360 (2020).
57. Jin, C., Song, C., Bjelland, J., Canright, G. & Wang, D. Emergence of scaling in complex substitutive systems. *Nature Human Behaviour* **3**, 837–846 (2019).
58. Leroi, A. M. *et al.* On revolutions. *Palgrave Communications* **6**, 1–11 (2020).
59. Tria, F., Loreto, V. & Servedio, V. D. P. Zipf's, heaps' and taylor's laws are determined by the expansion into the adjacent possible. *Entropy* **20**, 752 (2018).
60. Sinatra, R., Condorelli, D. & Latora, V. Networks of motifs from sequences of symbols. *Physical Review Letters* **105**, 178702 (2010).
61. Ha, L. Q., Hanna, P., Ming, J. & Smith, F. J. Extending zipf's law to n-grams for large corpora. *Artificial Intelligence Review* **32**, 101–113 (2009).
62. Ryland Williams, J. *et al.* Zipf's law holds for phrases, not words. *Scientific reports* **5**, 12209 (2015).
63. Csányi, G. & Szendrői, B. Structure of a large social network. *Physical Review E* **69**, 036131 (2004).
64. Glänzel, W. Characteristic scores and scales: A bibliometric analysis of subject characteristics based on long-term citation observation. *Journal of Informetrics* **1**, 92–102 (2007).
65. Milojević, S. Modes of collaboration in modern science: Beyond power laws and preferential attachment. *Journal of the american society for information science and technology* **61**, 1410–1423 (2010).
66. Milojević, S. Power law distributions in information science: Making the case for logarithmic binning. *Journal of the American Society for Information Science and Technology* **61**, 2417–2425 (2010).
67. Gerlach, M. & Altmann, E. G. Stochastic model for the vocabulary growth in natural languages. *Physical Review X* **3**, 021006 (2013).
68. Newman, M. E. & Watts, D. J. Scaling and percolation in the small-world network model. *Physical review E* **60**, 7332 (1999).
69. Watts, D. J. *Small worlds: the dynamics of networks between order and randomness*, vol. 36 (Princeton university press, 2004).
70. Gravino, P., Servedio, V. D., Barrat, A. & Loreto, V. Complex structures and semantics in free word association. *Advances in Complex Systems* **15**, 1250054 (2012).
71. Cattuto, C., Loreto, V. & Pietronero, L. Semiotic dynamics and collaborative tagging. *Proceedings of the National Academy of Sciences* **104**, 1461–1464 (2007).
72. Cattuto, C., Baldassarri, A., Servedio, V. D. P. & Loreto, V. Vocabulary growth in collaborative tagging systems. *arXiv preprint arXiv: 0704.3316* (2007).
73. Fink, T. & Reeves, M. How much can we influence the rate of innovation? *Science Advances* **5**, eaat6107 (2019).
74. Barabási, A.-L. & Albert, R. Emergence of scaling in random networks. *Science* **286**, 509–512 (1999).
75. Bianconi, G. & Barabási, A.-L. Competition and multiscaling in evolving networks. *Europhysics letters* **54**, 436 (2001).
76. Monechi, B., Gravino, P., Servedio, V. D. P., Tria, F. & Loreto, V. Significance and popularity in music production. *Royal Society Open Science* **4**, 170433 (2017).
77. Castaldo, M., Venturini, T., Frasca, P. & Gargiulo, F. Junk news bubbles modelling the rise and fall of attention in online arenas. *new media & society* **24**, 2027–2045 (2022).

78. Last.fm. Description page. <https://www.last.fm/about> (Accessed: January 2021).
79. Celma, O. *Music Recommendation and Discovery in the Long Tail* (Springer, 2010).
80. Last.fm. Music recommendation datasets for research. last.fm dataset - 1k users. <http://ocelma.net/MusicRecommendationDataset/lastfm-1K.html> (2009). Accessed: January 2021.
81. MusicBrainz. Home page. <https://musicbrainz.org/> (2022). Accessed: June 2022.
82. Hart, M. Project Gutenberg, 1971. *Project Gutenberg*, URL: <https://www.gutenberg.org> (1971).
83. Gerlach, M. & Font-Clos, F. A standardized project Gutenberg corpus for statistical analysis of natural language and quantitative linguistics. *Entropy* **22**, 126 (2020).
84. Porter, M. F. Snowball: A language for stemming algorithms (2001).
85. Porter, M. F. An algorithm for suffix stripping. *Program* (1980).
86. Paice, C. D. Another stemmer. *SIGIR Forum* **24**, 56–61 (1990).
87. Scholar, S. Home page. <https://www.semanticscholar.org/> (Accessed: January 2021).
88. Ammar, W. *et al.* Construction of the literature graph in semantic scholar. In *NAACL* (2018).
89. Scholar, S. Api license agreement. <https://api.semanticscholar.org/license/> (Accessed: June 2023).
90. Allen Institute for AI. Semantic scholar’s paper field of study classifier. [https://github.com/allenai/s2\\_fos](https://github.com/allenai/s2_fos). Accessed June 8, 2022 (2022).
91. Virtanen, P. *et al.* Scipy 1.0: fundamental algorithms for scientific computing in python. *Nature Methods* **17**, 261–272 (2020).
92. Tria, F., Crimaldi, I., Aletti, G. & Servedio, V. D. P. Taylor’s law in innovation processes. *Entropy* **22**, 573 (2020).
93. Pólya, G. Sur quelques points de la théorie des probabilités. In *Ann. Inst. Henri Poincaré*, vol. 1, 117–161 (1930).
94. Johnson, N. L. & Kotz, S. Urn models and their application; an approach to modern discrete probability theory (1977).
95. Mahmoud, H. Pólya urn models crc press. *Boca Raton FL* (2009).
96. Di Bona, G., Iacopini, I., Petralia, A. & Latora, V. Code used to download, process and analyse the data and the models at higher orders. <https://github.com/gabriele-di-bona/higher-order-heaps-laws> (2023).
97. King, T., Butcher, S. & Zalewski, L. *Apocrita - High Performance Computing Cluster for Queen Mary University of London* (2017).

## Acknowledgements

G.D.B. acknowledges support from the French Agence Nationale de la Recherche (ANR), under grant ANR-21-CE38-0020 (project ScientIA). A.P. and V.L. acknowledge support from the PNRR GRInS Project. I.I. acknowledges partial support from the James S. McDonnell Foundation 21<sup>st</sup> Century Science Initiative “Understanding Dynamic and Multi-scale Systems”. All computations have been performed via the High Performance Computing (HPC) cluster provided by Queen Mary University of London<sup>97</sup>. A.P. wishes to thank Roberto di Mari for several comments and suggestions.

## Author contributions

G.D.B., I.I., A.P., and V.L. designed the study. A.P. performed a preliminary investigation, collected in an early draft. G.D.B. carried out the data collection and performed the numerical simulations. G.D.B., A.B., and G.D.M. carried out the analytical calculations. G.D.B., A.B., I.I., and V.L. wrote the

manuscript. All authors contributed to analyze the data, discuss the results, define the proposed model, and revise the manuscript.

## Competing interests

The authors declare no competing interests.

## Additional information

**Supplementary Information** is attached to this manuscript.