

Supporting Information

Wasserman et al. 10.1073/pnas.1412198112

SI Text

Network. To control for biases in data reporting in IMDb (1), we choose to consider only the network of connections between films made in the United States. Furthermore, we only select connections designated on IMDb as references, spoofs, or features. We also limit our analysis to films released in 2011 or earlier. We obtain the IMDb film connections information—as well as information on country of production, primary language, and production companies—from plain text data files provided through ftp sites (2). We use a Python (version 2.7.8) program developed in-house to parse the relevant information from this file.

After entering the valid connections, we construct a network where each connection is a directed edge and each film is a node (Fig. 1). A link from movie A to movie B signifies that movie A cites movie B. To ensure proper maintenance of the timeline, we only include a connection from A to B if A was released in a later calendar year than B. As such, the network contains no links that are “forward” in time and no links between two films released in the same calendar year. Therefore, the resulting network is acyclic. We then take the largest weakly connected component of this network, known as the giant component, for our analysis.

Data. After the network is constructed, we count the number of times each film is cited (in-degree) and the number of citations each film makes (out-degree). We compute the PageRank value for each film in the network using the NetworkX Python package (version 1.8.1). The formula for the time lag of a citation is as follows:

$$t = y(k_{\text{out}}) - y(k_{\text{in}}), \quad [\text{S1}]$$

where $y(k)$ is the year of release of film k , and k_{out} and k_{in} are the films on the outgoing and incoming sides of an edge, respectively. After calculating the time lag for every edge in the network, we count the number of citations with time lag of at least 25 y that each film receives. This is the long-gap citation count.

We collect data on IMDb average user ratings and total numbers of votes for each film in the network through provided text files (2). Data on box office information and genre are also obtained through these files. We use Python programs developed in-house for parsing these files. The IMDb ID numbers for each film, which are necessary for accessing a film’s page on the IMDb website (www.imdb.com), are obtained through an in-house web scraping Python program using the BeautifulSoup package (version 4.3.2). We use this package in all our web-scraping processes.

We scrape Metacritic scores for films from web pages on the Metacritic website (www.metacritic.com). Each Metacritic web page is accessed via a film’s “critic reviews” page on the IMDb website, which contains a direct link to Metacritic if an aggregate review score exists for that film. We scrape Roger Ebert ratings for films from pages on Ebert’s official site (www.rogerebert.com). Each page on Ebert’s site is accessed through a film’s “external reviews” page on IMDb, which consists of user-added links to reviews of films on external websites. If Roger Ebert reviewed a film, a link to his review generally appears first on this page. We manually compile the list of films present in the National Film Registry (NFR) as the limited number makes this option possible (3).

Distribution Modeling. To generate null models for the distribution of time lags in the film connections network, we create Markov chain Monte Carlo simulations wherein the network undergoes

random rewiring (4–7) (Fig. S2). In each step of a simulation, two edges are selected at random from the network of films as candidates for rewiring. If the candidate connections “overlap”—that is, if at least one of the films of edge E was released in a calendar year in-between the years of release of the two films of edge F, noninclusive—then a swapping of connection nodes occurs. The “swapping” process consists of removing the chosen edges E and F from the network and replacing them with two new edges G and H, where edge G connects the outgoing film of edge E to the incoming film of edge F, and edge H connects the outgoing film of edge F to the incoming film of edge E. By allowing swapping between overlapping edges, we ensure that no back-in-time links are created. We forbid swapping if one of the edges created as a result of swapping already exists in the network. This process allows for random redistribution of edges while maintaining the in- and out-degrees of all of the nodes.

We use the simulation to generate the base null model—where the two randomly chosen edges are always swapped when it is legal to do so—as well as a null model with a bias toward shorter-length citations. In these latter simulations, a legal pair of randomly chosen edges undergoes swapping with probability q :

$$q = e^{\frac{\min(t_1, t_2) - \min(s_1, s_2)}{40}}, \quad [\text{S2}]$$

where t_1 and t_2 are the time lags of the two chosen edges and s_1 and s_2 are the time lags of the two edges if they were to be swapped. More specifically, if $t_1 = y_1 - z_1$ and $t_2 = y_2 - z_2$, where y_i and z_i are the years of the films connected by edge i , then $s_1 = y_1 - z_2$ and $s_2 = y_2 - z_1$.

In each run of a simulation, $20n_e$ iterations are performed—where n_e is the number of edges in the network. In total, we run 400 simulations, 200 with the base simulation and 200 with the biased simulation. We use Python programs developed in-house to run all rewiring simulations.

In addition, we use a theoretical formula for the time lag distribution of the unbiased null model, given nodes with specific in-degrees, out-degrees, and years of release:

$$\mathbb{E}(L_t) = \sum_{y \in \mathcal{Y}} \mathbb{E}(c_{y, y-t}), \quad [\text{S3}]$$

where L_t is the number of links with time lag t in the null model, \mathcal{Y} is the set of all years of release for films in the network, and $c_{y,z}$ is the number of links between films released in year y and films released in year z ($y > z$). The expected value of $c_{y,z}$ is determined by the following formula:

$$\mathbb{E}(c_{y,z}) = \frac{o_y i_z}{\sum_{j \in \mathcal{Y}} (i_j - o_j)} \prod_{k=z+1}^{y-1} \left(1 - \frac{o_k}{\sum_{j \in \mathcal{Y}} (i_j - o_j)} \right), \quad [\text{S4}]$$

where i_y and o_y are the sum totals of in-citations and out-citations, respectively, for films released in year y . We adapt this equation from Karrer and Newman’s formula for the expected number of edges between vertices in a directed acyclic graph with a fixed degree sequence (8).

Linear Regression. We narrow our focus to seven metrics: Roger Ebert rating, Metacritic score, IMDb average user rating, number of IMDb votes, total citation count, PageRank score, and long-gap

citation count. We calculate the adjusted R^2 values of linear regressions between each pair of considered measures using the statsmodels Python package (version 0.5.0). In our linear regression models, we use the base-10 logarithm of IMDb votes and PageRank score and the cube root of citations and long-gap citations. We opt to use cube root rather than log for citations and long-gap citations because many films have 0 values for these metrics, and positive values extend over several orders of magnitude. All regressions we perform in this paper apply these functions to these metrics.

Probit Regression. We perform probit regressions of the following form:

$$\text{inNFR} \sim \text{SigMetric}, \quad [\text{S5}]$$

where *inNFR* is the categorical variable representing whether or not a film is in the NFR (1 if it is in the NFR and 0 if it is not) and *SigMetric* is one of the seven metrics. For metrics with missing data—which are the expert-based metrics and the IMDb voting statistics—we apply the Heckman correction method (9, 10) to the probit regression, using R (version 3.0.2) and the sampleSelection package (version 1.0-2) (11). For metrics without missing data, we perform the regression with the statsmodels package in Python. We use probit instead of logit for this analysis because the sampleSelection package can only apply the Heckman correction for binary outcomes using probit.

When we apply the Heckman correction method, we use year of release and film genre as the dependent variables in the selection model equation. We note that genre is actually a set of 24 binary variables representing the 24 categorical film genres listed on IMDb. Films are not limited to being classified as one genre, and 11,661 of films in the network (or 75.6%) are categorized under two or more genres.

For all of these models apart from the long-gap citations model, we perform the regression on the subset of films released on or before 2003, as only films released on or before that year were eligible for nomination to the NFR in 2013 (3). For the long-gap citations model, we perform the regression on the subset of films made on or before 1986. The justification for the different subsets is that all films released after 1986 in our dataset have zero long-gap citations. (Our dataset only includes films released up to 2011, and the latest year that can possibly have a nonzero number of citations with a 25-or-more-year time lag is 1986.)

From the probit regression models, we obtain estimated probabilities for each film used to create the model. From these estimated probabilities, we assign a predicted value of 0 or 1 to each observation (0 if the probability is below 0.5, and 1 if the probability is greater than or equal to 0.5). We use the actual and predicted values to construct the classification table. We use the classification table to compute the balanced accuracy:

$$\text{Balanced Accuracy} = \frac{\text{Sensitivity} + \text{Specificity}}{2} = \frac{1}{2} \left(\frac{\text{TP}}{\text{TP} + \text{FN}} + \frac{\text{TN}}{\text{TN} + \text{FP}} \right), \quad [\text{S6}]$$

where TP, TN, FP, and FN are true positives, true negatives, false positives, and false negatives, respectively. We use the balanced accuracy instead of the true accuracy because the latter is strongly affected by the imbalance toward films not in the NFR versus those that are.

We also use the estimated probabilities to determine the receiver operating characteristic (ROC) curve (12). We calculate the area under the ROC curve (AUC) with the scikit-learn Python package (version 0.14.1). Finally, we use our probit regression results to calculate the pseudo- R^2 using Tjur's equation (13).

We obtain SD values for the balanced accuracy, AUC, and pseudo- R^2 of each metric using bootstrapping with 1,000 random samples. We use programs developed in-house in either Python or R to conduct all of the bootstrapping we do for this project.

Additionally, we repeat the same analysis detailed in this section, only instead of accounting for missing data, we ignore it and perform probit regression on the reported values (Table S1). This allows us to clearly present the effect of missing data and the Heckman correction.

Random Forest Classification. We perform Random Forest (RF) classification (14) using R and the randomForest package (version 4.6–10) (15). We conduct RF classification once with all seven aforementioned metrics as predictor variables, and another time with all metrics apart from long-gap citation count. In both cases, we use presence in the NFR as the binary response. We perform cross-validation by conducting 100 iterations of RF classification with each iteration using 80% of the data points, chosen randomly without replacement. We use the subset of films that were made in 1999 or earlier and have reported data for all seven metrics in our RF classification. This subset consists of 766 films. We conduct each classification iteration using 1,000 classification trees. From the cross-validated RF classification results, we obtain the mean and SD of variable importance—also known as the permutation importance—for each predictor.

Multivariate Regression. We perform two probit regressions using multiple independent variables. The first uses all metrics apart from long-gap citations as independent variables, whereas the second includes long-gap citations. In both regressions, the dependent variable is presence in the NFR. Also, both regressions are performed on the same subset of films used in RF classification. We evaluate the fit of the regression models by calculating the pseudo- R^2 with McFadden's equation (16). We perform these regressions with the statsmodels Python package.

We also repeat this same analysis but with logit instead of probit to demonstrate the minimal differences between the regression models (Table S2).

Citation Description Analysis. The brief notes that accompany some film connections on IMDb are not provided in the aforementioned plain text files, which we originally used to construct the connections network. Instead, we obtain these descriptions by scraping them from the actual IMDb movie connections pages. For each citation in the network, we check the cited film's connections page to see first whether the citing film is listed, and second whether a description is included with that citation. If a citing film is listed twice on the page and each listing has a description, then both descriptions are scraped. As with the initial construction of the network, we only scrape a description if the citation is classified as a reference, spoof, or feature.

After obtaining the citation descriptions, we proceed with two methods of analysis. In the first method, we take a small subset of highly cited films and, by hand, classify all of the annotations based on what they are citing. The subset of films we consider is the bottom 15 films from Table S3 (i.e., from *Bride of Frankenstein* to *Dirty Harry*). We classify annotated citations as “general” if the annotations merely refer to a film's title, title character, or plot, or if the citation is to numerous clips of the film. If an annotation is not general, then we classify the citation according to the part of the film to which it pertains, such as a specific scene, quotation, character, setting, or song. Two people independently classified the annotated citations for these films. The two people differed by no more than two citations in any classification for any film. The results of this manual classification are shown in Table S4.

In the second method, we use the `token_set_ratio` function from the `fuzzywuzzy` Python package (version 0.3.2) to perform

comparisons between citation descriptions. Initially, we clean all of the descriptions by removing all punctuation—apart from hyphens and apostrophes in-between letters—and converting all alphabet characters to lowercase. For each film with a minimum of 20 annotated citations, we compare every pair of descriptions using the `token_set_ratio` function, which returns an integer value indicating the similarity of two strings, with 0 being the least similar and 100 being the most similar. Thus, for a film with c annotated citations, we obtain $\binom{c}{2}$ similarity values. Taking the average of all of the similarity values for a film gives us the “mean similarity” for that film’s citation descriptions.

To compensate for differing numbers of descriptions and varying lengths of strings, we perform bootstrapping wherein all of the words in all of the descriptions for a specific film are randomly redistributed while keeping the number of words in each description constant. We then perform the aforementioned process for computing the mean similarity on the jumbled citation descriptions. We perform 500 randomization iterations for each film with a minimum of 20 annotated citations. We then obtain a mean and SD for all of the randomized mean similarities for a film, as well as a Z score for the mean similarity of the actual descriptions. We perform linear regressions comparing the Z scores to the results of manual classification.

1. Wasserman M, et al. (2014) Correlations between user voting data, budget, and box office for films in the Internet Movie Database. *J Am Soc Inf Sci Technol*, 10.1002/asi.23213.
2. Internet Movie Database (2012) Alternative Interfaces. Available at www.imdb.com/interfaces. Accessed October 26, 2012.
3. Library of Congress (2014) National Film Registry. Available at www.loc.gov/film/filmnfr.html. Accessed April 11, 2014.
4. Maslov S, Sneppen K (2002) Specificity and stability in topology of protein networks. *Science* 296(5569):910–913.
5. Milo R, et al. (2002) Network motifs: Simple building blocks of complex networks. *Science* 298(5594):824–827.
6. Milo R, Kashtan N, Itzkovitz S, Newman MEJ, Alon U (2004) On the uniform generation of random graphs with prescribed degree sequences. [arXiv:cond-mat/0312028](http://arxiv.org/abs/cond-mat/0312028).
7. Carstens C (2013) Motifs in directed acyclic networks. *2013 International Conference on Signal-Image Technology and Internet-Based Systems* (IEEE Computer Society, Los Alamitos, CA), pp 605–611.
8. Karrer B, Newman MEJ (2009) Random graph models for directed acyclic networks. *Phys Rev E Stat Nonlin Soft Matter Phys* 80(4 Pt 2):046110.
9. Heckman JJ (1976) The common structure of statistical models of truncation, sample selection and limited dependent variables and a simple estimator for such models. *Ann Econ Soc Meas* 5(4):475–492.
10. Heckman JJ (1979) Sample selection bias as a specification error. *Econometrica* 47(1): 153–161.
11. Toomet O, Henningsen A (2008) Sample selection models in R: Package sample-Selection. *J Stat Softw* 27(7):1–23.
12. Zweig MH, Campbell G (1993) Receiver-operating characteristic (ROC) plots: A fundamental evaluation tool in clinical medicine. *Clin Chem* 39(4):561–577.
13. Tjur T (2009) Coefficients of determination in logistic regression models—a new proposal: The coefficient of discrimination. *Am Stat* 63(4):366–372.
14. Breiman L (2001) Random forests. *Mach Learn* 45(1):5–32.
15. Liaw A, Wiener M (2002) Classification and regression by randomForest. *R News* 2(3): 18–22.
16. McFadden D (1974) *Conditional Logit Analysis of Qualitative Choice Behavior*. *Frontiers in Econometrics*, ed Zarembka P (Academic, New York), pp 105–142.

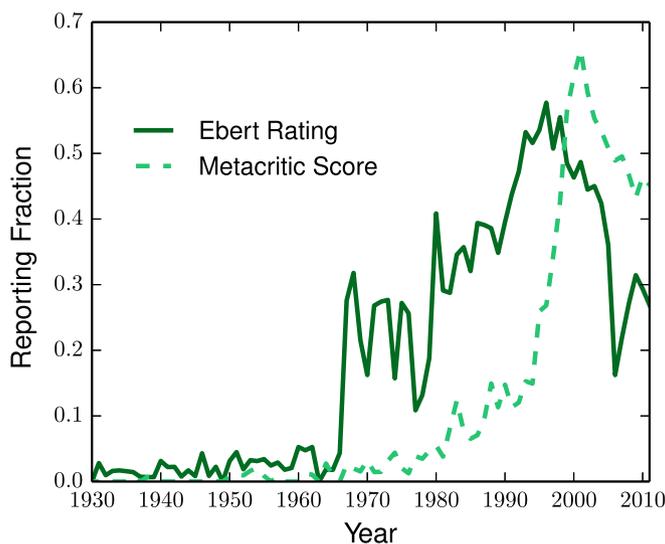


Fig. 51. Fraction of reported critic data values by year.

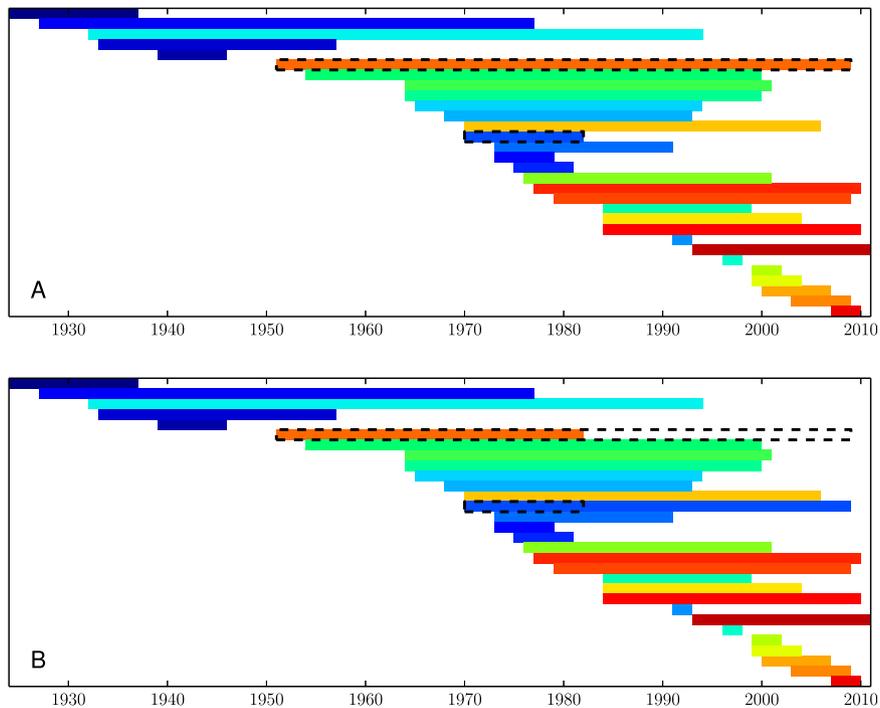


Fig. S2. One step of rewiring simulation. The diagram depicts a randomly chosen sample of edges from the directed network of film connections. Each bar represents an edge in the sample. The right end of each bar indicates the year of release for the citing film. The left end of each bar indicates the year of release for the cited film. (A) Two edges are selected at random as candidates for swapping. (B) If the two chosen edges overlap, they are removed from the network and replaced with two new edges that connect the outgoing film of one original edge to the incoming film of the other original edge. The black dotted lines represent the originally chosen candidate edges, now removed from the network.

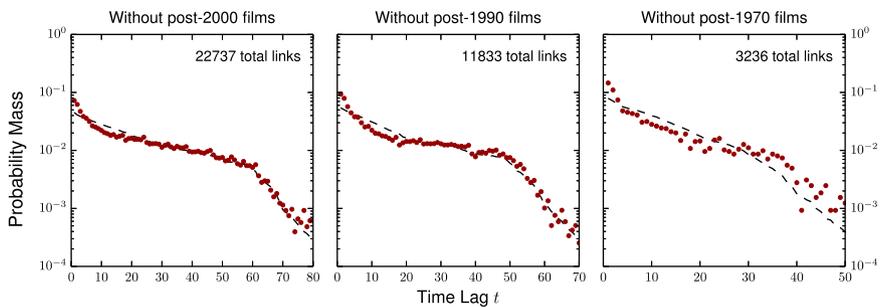


Fig. S3. Distribution of time lag with exclusions. Probability mass function of the time lag of connections in the film connections network, discounting all films made after 2000, after 1990, and after 1970. Brown points represent the actual distributions. Dashed black lines represent the unbiased null model distribution, calculated with Eq. S3.

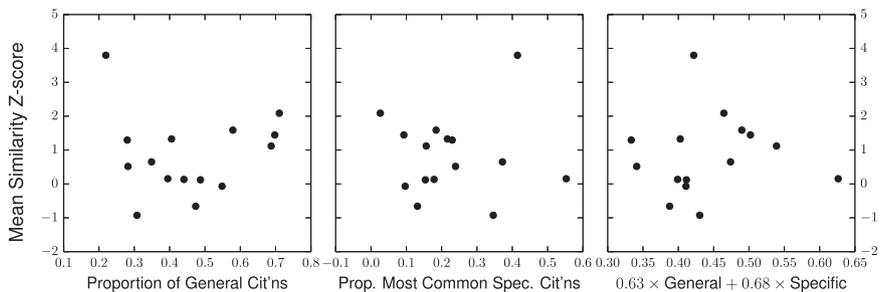


Fig. S4. Mean similarity Z score versus annotated citation classifications. Scatterplots comparing the mean similarity Z scores of citation annotations to the fractions of annotations under certain classifications for 15 highly cited films (See Table S4). The fractions considered are the proportion of general citations (Left), the proportion of most common specific citations (Center), and the best-fitting linear combination of the two proportions obtained through ordinary least-squares regression (Right). No adjusted R^2 value is positive.

Table S3. Films with most long-gap citations

Title	Year	LGC*	NFR year [†]
<i>The Wizard of Oz</i>	1939	565	1989
<i>Star Wars</i>	1977	297	1989
<i>Psycho</i>	1960	241	1992
<i>Casablanca</i>	1942	212	1989
<i>Gone with the Wind</i>	1939	198	1989
<i>King Kong</i>	1933	191	1991
<i>Frankenstein</i>	1931	170	1991
<i>The Godfather</i>	1972	162	1990
<i>Citizen Kane</i>	1941	143	1989
<i>2001: A Space Odyssey</i>	1968	143	1991
<i>Jaws</i>	1975	129	2001
<i>Night of the Living Dead</i>	1968	122	1999
<i>It's a Wonderful Life</i>	1946	109	1990
<i>The Graduate</i>	1967	97	1996
<i>Vertigo</i>	1958	92	1989
<i>Snow White and the Seven Dwarfs</i>	1937	91	1989
<i>Dr. Strangelove or: How I Learned to Stop Worrying and Love the Bomb</i>	1964	91	1989
<i>Dracula</i>	1931	90	2000
<i>The Maltese Falcon</i>	1941	80	1989
<i>Bambi</i>	1942	79	2011
<i>The Exorcist</i>	1973	78	2010
<i>Taxi Driver</i>	1976	71	1994
<i>Sunset Blvd.</i>	1950	70	1989
<i>Planet of the Apes</i>	1968	69	2001
<i>Deliverance</i>	1972	66	2008
<i>The Sound of Music</i>	1965	61	2001
<i>Bride of Frankenstein</i>	1935	58	1998
<i>Singin' in the Rain</i>	1952	57	1989
<i>Apocalypse Now</i>	1979	57	2000
<i>The Texas Chain Saw Massacre</i>	1974	57	
<i>Rebel Without a Cause</i>	1955	57	1990
<i>Star Wars: Episode V—The Empire Strikes Back</i>	1980	56	2010
<i>North by Northwest</i>	1959	54	1995
<i>Rear Window</i>	1954	54	1997
<i>Mary Poppins</i>	1964	54	2013
<i>Pinocchio</i>	1940	53	1994
<i>Willy Wonka & the Chocolate Factory</i>	1971	52	
<i>The Seven Year Itch</i>	1955	51	
<i>Rosemary's Baby</i>	1968	51	
<i>West Side Story</i>	1961	51	1997
<i>Dirty Harry</i>	1971	51	2012

*Long-gap citation count.

[†]Year inducted into the NFR (3).

