# Why Name Popularity is a Good Test of Historicity[1*]

## A Goodness-of-fit Test Analysis on Names in the Gospels and Acts

Luuk van de Weghe, Ph.D. ORCHID: 0000-0001-8710-503X

Independent scholar, Port Angeles, WA, 98362, USA

luukvandeweghe@gmail.com


Jason Wilson, Ph.D., ORCID: 0009-0003-7734-1110

Department of Math and Computer Science, Biola University, La Mirada, CA, 90639, USA

jason.wilson@biola.edu

## Abstract

Are name statistics in the Gospels and Acts a good test of historicity? Kamil Gregor and Brian Blais, in a recent article in *The Journal for the Study of the Historical Jesus*, argue that the sample of name occurrences in the Gospels and Acts is too small to be determinative and that several statistical anomalies weigh against a positive verdict. Unfortunately, their conclusions result directly from improper testing and questionable data selection. Chi-squared goodness-of-fit testing establishes that name occurrences in the Gospels and Acts fit into their historical context at least as good as those in the works of Josephus. Additionally, they fit better than occurrences derived from ancient fictional sources and occurrences from modern, well-researched historical novels.

## Keywords

authenticity, Bauckham, Gospels and Acts, Gregor and Blais, onomastics, statistics, goodness-of-fit

Our topic has its roots in the publication of Tal Ilan's 2001 *Lexicon of Jewish Names in Late Antiquity (Part 1: Palestine 330 BCE – 200 CE).*[2] This lexicon, like a telephone book, catalogs the names of approximately 2500 persons into a single volume, giving us data related to name origin and name popularity statistics of Palestinian Jews living around the time of Jesus. For this enormous undertaking, and her three subsequent volumes, we owe her a debt of gratitude.[3]

New Testament scholar Richard Bauckham soon took advantage of this database to aid in his own radically innovative study of names in the Gospels and Acts (GA). He noticed how names were conserved but sometimes forgotten throughout the redactional stages of the Synoptic traditions, how they emerged in varied forms within parallel lists, and how they clustered around certain POV (point-of-view) perspectives.[4] Bauckham also observed, using Ilan's lexicon, that the name popularity statistics in GA generally conform to the population statistics attested to in Ilan's database.[5]  Figure 1 below portrays Bauckham's insight in a graphical form, with the tops of the bars representing the percentages from Ilan's database. Although not perfect,[6] the fit is striking and was the genesis of his hypothesis.

---

[2] Tal Ilan, *Lexicon of Jewish Names in Late Antiquity. Part I, Palestine 330 BCE–200 CE* (Tübingen: Mohr Siebeck, 2002).

[3] Her three other volumes are: Tal Ilan, *Lexicon of Jewish Names in Late Antiquity. Part III, The Western Diaspora 330 BCE–650 CE* (Tübingen: Mohr Siebeck, 2008); Tal Ilan, *Lexicon of Jewish Names in Late Antiquity. Part IV, The Eastern Diaspora 330 BCE–650 CE* (Tübingen: Mohr Siebeck, 2011); Tal Ilan, *Lexicon of Jewish Names in Late Antiquity. Part II, Palestine 200–650 CE (Tübingen:* Mohr Siebeck, 2012).

[4] Richard Bauckham, *Jesus and the Eyewitnesses: The Gospels as Eyewitness Testimony* (Grand Rapids: Eerdmans, 2017, 2nd ed.), pp. 42, 44, 46-55, 156–64, etc. Another recent example is Bauckham, 'Eyewitnesses and Healing Miracles in the Gospel of Mark,' *The Biblical Annals* 10 (2020), pp. 341-354, esp. p. 351.

[5] Bauckham, *Jesus and the Eyewitnesses*, pp. 67–84.

[6] Eleazar (Lazarus) and Jacob (James) are off, although they fall within the confidence intervals, depending on how many name occurrences are considered in the Gospels and Acts sample. The data for Figure 1 is discussed in Section 2 below. Following Bauckham, this graph includes occurrences from Josephus and GA within the Ilan data,
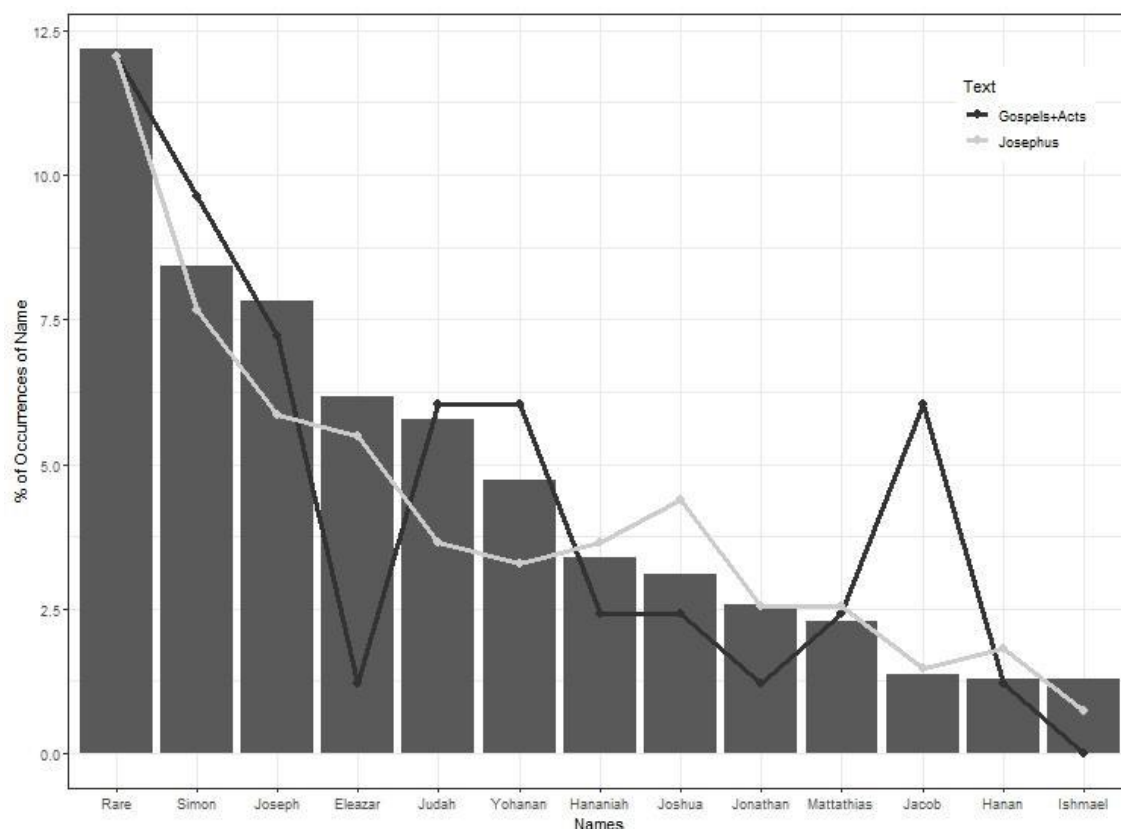
*Figure 1. Top 12 names %'s from Ilan-1 shown for Gospels+Acts and Josephus. The bar heights represent the height of the full Ilan-1 % of the name frequencies of the 12 most frequently occurring names. The %'s for Gospels+Acts, and Josephus, are superimposed with a dot. A connective line was added for ease of viewing. In addition, the % of rare names (=1 occurrence) was added.*

A question recently posed by Kamil Gregor and Brian Blais (GB) was whether this insight could withstand formal statistical analysis. They call the significance of Bauckham's observations into question while simultaneously furthering the conversation in this novel area of study.[7]

Few studies have otherwise sought to apply statistical analysis to onomastic data, and even here the aims and subject matter were quite dissimilar from Bauckham's project.[8] These

---

but in our statistical analysis we exclude them (see Section 3.1 below, footnote 39). We briefly revisit this graph in Section 2.

[7] Kamil Gregor and Brian Blais, 'Is Name Popularity a Good Test of Historicity?,' *JSHJ* 21.3 (2023), pp. 171–202.

[8] See Dan Zao, 'Snack Names in China: Patterns, Types, and Preferences,' *Names* 69.4 (2021), pp. 13-20 (p. 16); Alessandra Minello, Gianpiero Dalla-Zuanna, and Guido Alfani, 'The Growing Number of Given Names as a Clue to the Beginning of the Demographic Transition in Europe,' *Demographic Research* 45 (2021), pp. 187–220 (p. 205); M.

parallel studies also lack in-depth discussions on statistical methodology. The scarcity of this literature highlights the extent to which Bauckham's approach – which had the potential of combining the hard and soft sciences in a single endeavor – was innovative. GB advance the conversation in several ways. First, they draw our attention to Ilan's complete 4-volume database in its current machine-readable form online.[9] Second, they rightly lament the lack of statistical rigor in Bauckham's analysis.[10] Finally, they point out several significant miscalculations by Bauckham and do well to constrain their analysis to Palestinian Jews from 4 BCE–73 CE to capture an onomastic snapshot of the time of Jesus.[11]

Critical missteps, however, undermine GB's assessment of names in GA. Due to their isolated focus on contested names, their improper testing methodology, and their lack of adequate data pre-processing, they fall short of providing the "suitably robust" statistical analysis that their piece claims.[12] After presenting our statistical analysis, we provide an in-

---

Depauw and W. Clarysse, 'How Christian was Fourth Century Egypt? Onomastic Perspectives on Conversion,' *Vigiliae Christianae* 67.4 (2013), pp. 407–35 (p. 424).

[9] For the full machine-readable database, see: https://github.com/hlapin/eRabbinica/tree/master/ilanNames (accessed January 9, 2024).

[10] Gregor and Blais, 'Name Popularity,' p. 180; cf. Bauckham, *Jesus and the Eyewitnesses*, pp. 67–84.

[11] The period of 330 BCE–200 CE from Ilan, *Part I*, needed to be narrowed to achieve a "snapshot impression of an onomastic situation," as Ilan calls it (Ilan, *Part I*, p. 50). See also, Luuk van de Weghe, 'Name Recall in the Synoptic Gospels', *NTS* 69.1 (2023), pp. 95–109 (p. 96); the careful reader is immediately struck by the 371 name occurrences noted in this piece versus the 2181 noted by GB. The main reason for this discrepancy is not articulated in 'Name Popularity' and appears to be the difference between an inclusive or exclusive approach in dating methodology. Take, for example, GB's time period of 4 BCE–73 CE. Many of Ilan's entries are dated broadly (e.g., –135 CE, –200 CE, –70 CE, etc.). Few are dated specifically (29 CE, 63 CE, etc.). In Van de Weghe, 'Name Recall,' an exclusive approach is used, only counting a name occurrence if it could be dated specifically within 30 BCE–90 CE (pp. 105-06; note that p. 105 contains a typographical error and should read 371 names, not 391, per the table on p. 106). Gregor and Blais apparently take an inclusive approach, counting every name occurrence that cannot specifically be shown *not* to occur within their time window. In this article, we follow GB's timeframe and methodology.

[12] 'Name Popularity,' p. 202. One immediate concern is the lack of definition around "historicity." GB do not articulate but seem to imply a narrow definition, as if "historicity" equals the substantiation of every contested

depth critique of GB in Section 4 of this paper. Their study also contains several highly contentious speculations, but because our analysis broadly undermines the premises of almost every one of these there will be no need to address them in-depth.[13]

Below, we demonstrate that the Gospels and Acts, alongside Josephus, reflect a population of historical Palestinian Jewish names from 4 BCE–73 CE as opposed to a population of Diaspora males from the same period. Among other conclusions drawn from our analysis, we also demonstrate that name samples from the most robust historical novels available to us – novels that did indeed rely on Josephus, GA, and the Hebrew Bible for their naming practices – were not able to achieve naming statistics that significantly reflect Palestinian Jews from 4 BCE–73 CE. This latter analysis implicitly responds to a historically nuanced scenario envisioned by GB wherein the Gospel authors invented certain names but also relied on historical source material and population distributions found in Josephus.[14] In short, we show that name statistics in GA fit into their historical context well, and that they fit significantly better than fictitious samples.

## 1. Data

We incorporate seven sets of data into our analysis: two sets are reference distributions; five are test distributions. Reference distributions represent the real populations that our test distributions are compared against. For example, we will analyze name occurrences in GA (test

---

name. Bauckham's view, as he admitted, never aimed for such a high bar (Bauckham, *Jesus and the Eyewitnesses*, p. 544). A mere statistical analysis alone could not prove that every name in GA is historical, and a statistical analysis is not needed to disprove this point. For our purpose, we suggest treating "historicity" on par with "authenticity," and then in a broad sense also suggested by Bauckham when he writes, "all the evidence indicates the general authenticity of the personal names in the Gospels" (Bauckham, *Jesus and the Eyewitnesses*, p. 84). A widespread fitness of names would naturally lean in favor of the historicity of each individual name (see Van de Weghe, 'Name Recall,' pp. 108-09; Section 3.2 below), but it is not our aim here to defend so strong a position.

[13] These speculations include, but are not limited to, Dennis MacDonald's doublets, the relationship of the onomastic data of Luke-Acts to the writings of Josephus, and the alleged greater propensity of fictionalizers to invent popular versus rare names.

[14] Gregor and Blais, 'Name Popularity,' p. 188.

distribution) and Josephus (test distribution) against name occurrences from two populations (reference distributions) in Ilan's database. The purpose is to see how a given test distribution fits, or does not fit, into a reference distribution.

### 1.1 Reference Distributions

We accept the Ilan reference distribution for Palestinian Jewish males (4 BCE–73 CE) that GB scraped from Hayim Lapin's machine-readable database. We independently scraped this data, but due to minor variations between our scraped dataset and GB's,[15] as well as their valuable inclusion of the Josephus data, we chose to use the GB dataset for Ilan Vol. 1 (Ilan-1) with slight modifications, resulting in 2185 name occurrences.[16] As a brief point of clarification, a name occurrence is the number of times a name is attached to a unique person. For example, Simon has eight occurrences in GA because it is attached to eight persons. But Peter also has one occurrence, even though that name is attached to someone already named Simon. This must be distinguished from the number of times a person's name is mentioned in a text, which is irrelevant to our exercise. A title like "Simon Peter, also called Cephas" would give three occurrences, even though it only refers to one person and regardless of how many times this person might be referenced in our source material. Finally, we only consider males in our

---

[15] GB has 2181 names. They described three different sets of names they excluded from the complete set of names from Ilan, *Part 1*: Nicknames, Fictitious names, and Other. Following GB, there were 58 Nicknames in Vol. 1 that occur only once in the entire Ilan database. This brought the list to 2488 names. The "Other" exclusions included non-Palestinian people, such as those from Tarsus, and non-Jews, such as proselytes. Excluding all the Other and the Fictitious brought the number down to 2069 names. This is 112 less than GB. Since there were over 40 different categories in Other, some of which were a judgment call on whether they fit the non-Palestinian and non-Jew category, it is likely that the source of the discrepancy between the two datasets lies here. Such variation is not uncommon when scraping. Understandably, GB did not provide their list of exclusion categories, so it is neither possible nor worthwhile to further attempt to unify the two databases and we opted to follow GB.

[16] We added Qaifa (Caiaphas), which should have been included; although technically a family name, it evidently functions as a personal name in the relevant contexts. We also added the six deacons and removed Alexander, Rufus, and Bartimaeus for reasons discussed below. This resulted in 2185 name occurrences (2181, per GB, plus Qaifa, plus six deacons, minus Alexander, Rufus, and Bartimaeus).

analysis for the mere reason that antiquity's androcentric focus limits the available data for female name occurrences.[17]

A second reference distribution covers the names of Western Diaspora Jewish males from Tal Ilan Vol. 3 (Ilan-3) within the same time parameters as GB's Ilan-1 (4 BCE–73 CE). We choose Ilan-3 as an alternative reference distribution to further situate Bauckham's claim that someone writing in the Diaspora would be exposed to significantly different Jewish name distributions than those of Palestine. The data from this reference distribution justifies Bauckham's claim, showing a significantly high frequency of Greek versus biblical names (discussed further in Section 3.1). Among the top ten names, only the name Simon from Ilan-1 is present. The top ten names are, in order, Shabtai, Dositheus, Ptolemaius, Alexander, Simon, Gaius, Theodotus, Julius, Theodorus, and Philon (cf. Figure 1).[18]

### 1.2 Test Distributions

Aside from two reference distributions, we consider five test distributions. The first is the Gospels and Acts (GA). Like GB, we follow Bauckham's dataset, with three minor adjustments. Following Ilan, and in agreement with GB, we include the six deacons from Acts 6:5 (Stephen, Philip, Procorus, Nicanor, Timon, and Parmenas).[19]  Additionally, we exclude Rufus and Alexander (Mark 15:21) because, being the sons of Simone of Cyrene, they should likely not be considered Palestinian. Lastly, we remove Bartimaeus from GA due to inconsistencies associated with leaving it in.[20] This results in 82 name occurrences for GA (79, per Bauckham,

---

[17] These clarifications are also discussed by GB ('Name Popularity,' pp. 174-75) and in Van de Weghe, 'Name Recall,' p. 99.

[18] All data and *R* code used for the analysis, with some light explanation, is included in the Data and Supplementary Materials. https://doi.org/10.7910/DVN/N7API1

[19] GB state that there is no reason to assume the deacons, aside from Nicolas, are from the Diaspora ('Name Popularity,' p. 185).

[20] We should not consider Bartimaeus as a name in GA because we already consider Timaeus as a name, which is how Ilan designates all the Bar-names – that is, by their patronyms. While Mark 10:45 contains the clause, "Bartimaeus (which means 'son of Timaeus')," this simply clarifies what is implicit in all the Bar-names. In all other cases, Bauckham agrees with Ilan's designations, and while Bartimaeus appears to function as a name in Mark

plus six, minus three). A second test distribution we consider are the name occurrences in Josephus; we follow the helpful sample scraped by GB, with the addition of Qaifa.

Our third test distribution is the complete set of fictitious occurrences found in Ilan Vol. 1 (Ilan-1F). This sample derives from all name occurrences from 4BCE–73 CE belonging to Palestinian Jewish males that Ilan deemed fictitious. This pool springs from apocryphal gospel material as well as rabbinic material only loosely tied to authentic traditions.[21] Unlike the hypothetical uniform distribution considered by GB (discussed below), this sample derives from concrete data. We want to test how well the statistics in Ilan-1F fit, or do not fit, into the Ilan-1 reference distribution. This will give us an impression of how ancient fictionalizers performed in attempting to create authenticity in their naming practices for Palestinian Jewish males around the time of Jesus. Since this sample lumps many smaller discrete samples together, it allows data from fictitious sources with scarce onomastic data to be considered.

Our fourth test distribution is a complete sample of Palestinian Jewish male name occurrences from the novels *Ben Hur* and De Wohl's *The Spear*. Here we want to determine how well-researched historical novels perform in acquiring onomastic verisimilitude. This is significant, because GB imagined a scenario wherein GA could achieve an appropriate reflection of Ilan-1 by inventing names inspired by the works of Josephus as well as incorporating actual names of historical persons from various traditional materials.[22]

---

10:45, Bartholomew, which both Ilan and Bauckham designate by the patronym "Ptolemy," functions the same way (Mark 3:18). If we catalogue this name under the form "Bartimaeus," we should also catalogue Barsabbas, Bartholomew, etc., under the forms with their Bar- elements included rather than under the patronyms they derive from, but doing so would render all these names rare and go contrary to the practice followed by Ilan everywhere and by Bauckham everywhere but in this instance.

[21] It also includes occurrences from several inscriptions and characters from stories such as Jesus' Lazarus from Luke 16:20. Regarding post-Talmudic literature, Ilan writes, "some of these stories may be based on authentic traditions, but with the passage of time and the literary nature of these compositions this is not very likely" (*Part I*, p. 48).

[22] Gregor and Blais, 'Name Popularity,' pp. 188-89.

Louis De Wohl's *The Spear* is a meticulously researched historical novel set around the time of Jesus, incorporating many historical and fictional Palestinian Jewish characters.[23] Lewis Wallace's 1880 *Ben Hur*, while containing fewer Jewish names than *The Spear*, is a more familiar alternative. We therefore analyze all male name occurrences for both. Each work shows a wealth of historical data and research as well as unambiguous signs of dependence on Josephus and on the Gospels and Acts. Over 35 percent of their name occurrences appear to be directly influenced by these two sources.[24] Regarding Second Temple Judaism, Lew Wallace, the author of *Ben Hur*, claims to have visited the Library of Congress in 1873, researching, "everything on the shelves relating to the Jews."[25] These historical novels provide a good test case for what a well-informed inventor of Palestinian names might achieve in terms of appropriate naming patterns.

Finally, we consider a test distribution consisting of a uniform sample of 52 occurrences. GB suggested an additional scenario of an inventor randomly selecting names from a uniform distribution of the 454 Palestinian names that made up their Ilan-1 reference distribution; they argued that their sample of contested GA occurrences (53) was too small to significantly distinguish it from a random sample from Ilan-1.[26] We want to determine if this claim is true using our methodology, even though we take issue with their scenario (see Section 4.3).

### 1.3 Data Pre-Processing

Unfortunately, GB ignore data from names that have low occurrences in the GA test distribution: what they call "white noise."[27] This common problem in data analysis is easily overcome by the widely used pre-processing procedure called data binning. Data binning takes

---

[23] We thank Lydia McGrew for bringing De Wohl's *The Spear* to our attention for this analysis.

[24] 32 out of 86 occurrences are persons specifically mentioned in GA or Josephus: 37.31%. In Ben Hur alone, the percentage is even higher (19 out of 32 occurrences: 59.38%).

[25] Lew Wallace, *Lew Wallace: An Autobiography*, Vol. 2 (Harper & Brothers: London, 1906), p. 891.

[26] The discrepancy of our 52 names vs. their 53 names results from the removal of Bartimaeus from our distributions. For their discussion of this scenario, see Gregor and Blais, 'Name Popularity,' pp. 191–95.

[27] Gregor and Blais, 'Name Popularity,' p. 191; also, see our comments in Section 4 below.

data from an intricate set and categorizes it into discrete bins for the purpose of simplification, noise reduction, and enhanced analysis and modeling.[28] This technique allows us to lump even low-frequency occurrences from test distributions into categories.

Consider, for example, the data observed in the beginning of this paper in Figure 1. While Figure 1 gives great details of the most popular and the rarest name occurrences, it provides zero detail about more than half of the names in GA; had we not binned the rare names for this graph, it would represent less than 25% of names in GA. One might compare this to performing an in-depth analysis on voting preferences among millennials in order to draw conclusions about the voting trends in all demographics. Clearly, focusing only on that group will provide more detail about how particular ages are voting, but it is not the proper test for making decisions regarding the demographics of an entire population.

Especially with smaller sample sizes, as those we are dealing with, some mismatches are inevitable (e.g., Eleazar and Jacob in Figure 1). But Gregor and Blais claim that "even if there were only one Simon in Gospels-Acts, this would not fit the distribution of name popularity in the contemporary population statistically significantly worse than what we actually observe in Gospels-Acts! The same would be true even if *every* name among contested Gospels-Acts characters appeared only once."[29] But that is only true under their model, which insufficiently weighs the cumulative force that each low occurrence would have within a model like the one we use in this paper. Binning also allows us to consider more subtle features of population distribution than merely the most popular or most unusual names (See further discussion below, Section 3.2).

For our data pre-processing, we considered the following features: objectivity, consistency, and thoroughness. Regarding objectivity, we binned all data *relative to each reference distribution*; it was the reference distribution rather than the test distribution that determined the sizes and numbers of bins. This generated consistency in how we analyzed each test distribution. We utilized a method called equal frequency binning, meaning in our case that

---

[28] It is standard statistical practice. See a formal discussion of this procedure in the reference text by Alan Agresti, *Categorical Data Analysis* (Hoboken, NJ: Wiley & Sons, 2002), pp. 174-177.

[29] Gregor and Blais, 'Name Popularity,' p. 194.

we binned name frequencies into sets with approximately the same amount of name occurrences from our reference distributions.[30] Because our smallest possible bin (1 occurrence) encompassed approximately 1/6 of occurrences, the result was six bins.[31] This generated consistency across our tests of the popularity data. Our procedure will become clearer in Section 3, wherein we model our data.

Additionally, we improved on GB's analysis by considering name origin statistics in addition to name frequency statistics. Tal Ilan categorized name occurrences into eight categories according to name origin: Biblical, Greek, Latin, Persian, Egyptian, Arabian, Semitic-Hebrew, and Semitic-Greek. A strong case could be made for combining the latter two categories into a single category (Semitic), but because it would not significantly impact our analysis we opted to follow Ilan's designations. We analyze these categories as an independent exercise from our name frequency analysis.

## 2. Method

### 2.1 Goodness-of-fit Tests

For the analysis, we chose to use the most common and widely used statistical test of fit between a categorical test distribution and a categorical reference distribution: the chi-squared goodness-of-fit test.[32] For example, suppose that someone wanted to determine whether a die

---

[30] This was accomplished by taking the minimum root mean square error for different distributions of the name frequencies and selecting the one with the smallest root mean square error.  Because of the uneven name frequencies, the result may not be as even as one might expect (see Figure 2).  For more on equal frequency binning, see the seminal paper on this topic, H. B. Mann and A. Wald, 'On the Choice of the Number of Class Intervals in the Application of the Chi Square Test,' *Ann. Math. Stat*. 13.3 (1942), pp. 306–17.

[31] It was a little less for Ilan-1 and a little more for Ilan-3.  Another reason for six bins was that it was about right for satisfying the conditions of the chi-square goodness of fit test in all tests: no case fewer than one expected observation and fewer than 20% of cases have below five expected observations.

[32] Alan Agresti, *Categorical Data Analysis* (Hoboken, NJ: Wiley & Sons, 2002), pp. 22–26. One alternative considered was a chi-square test of independence, which would treat both the test and reference distributions as samples (the goodness-of-fit test treats the reference distribution as the ground truth). For thoroughness, we ran the chi-square tests of independence. The p-values were quite close to the goodness-of-fit p-values shown in Table

was fair, or not. The reference distribution would reflect a probability of 1/6 (i.e. "fair") for rolling a 1, 2, 3, 4, 5, or 6.  Suppose we rolled the die 60 times and obtained the observed sample in Table 1.[33]

| Face | 1 | 2 | 3 | 4 | 5 | 6 | | Total |
|---|---|---|---|---|---|---|---|---|
| Reference probability | 1/6 | 1/6 | 1/6 | 1/6 | 1/6 | 1/6 | | 6/6=1 |
| Observed | 5 | 8 | 9 | 8 | 10 | 20 | | 60 |
| Expected | 10 | 10 | 10 | 10 | 10 | 10 | | 60 |
| Difference | -5 | -2 | -1 | -2 | 0 | 10 | | 0 |

Table 1. Example of testing whether a 6-sided die is fair.

Based on the observed values, do you think the die is fair, or not? The test works by using the reference distribution to calculate the expected value, which in this case is (1/6)60 = 10 rolls expected for each of 1, 2, 3, 4, 5, and 6. The observed values are compared against the expected values and the probability value (*p-value*) is calculated.  The *p-value* is defined as the probability that the differences between the observed and expected values would be as great as the differences we observed, or greater, if the data truly came from the reference distribution. In this case, the *p-value = 0.0199*, meaning that there is about a *2%* probability that we would observe differences this large, particularly five 1's and twenty 6's, if the die is truly fair. In other words, this means that there is about a 2% chance of concluding the die is not fair, when actually it is.

---

3. They were usually slightly larger than those of the goodness-of-fit, but not always, and in no cases had any substantial divergence which would alter our conclusions.  They are shown in the Supplementary Materials.

[33] This particular example is taken from Wikipedia's Pearson's chi-squared test page, https://en.wikipedia.org/wiki/Pearson%27s_chi-squared_test. Accessed January 12, 2024.

Does this prove that the die is fair, or not fair? No, it only offers a probabilistic answer. In current scientific practice, p-values are compared against an arbitrary pre-selected benchmark, often *0.05* or *5%*. If the p-value falls below the benchmark, it is concluded that the sample distribution does not come from the reference distribution and the test is said to be **statistically significant.**[34] In this example we would have observed *p-value = 0.02 < 0.05*, meaning we conclude the die is not fair. The benchmark is the largest probability of making a false positive error that the researcher is willing to accept. Accompanying the statistical conclusions of any properly conducted statistical study should be all additional available scientific evidence and reasoning which are combined using logic to draw final conclusions. That is the process we will use in this paper.

There is one additional remark needed regarding the benchmark against which we will compare our *p-values*. When multiple tests are performed, as in our case, the so-called Multiple Testing Problem occurs. This problem is best captured in a cartoon[35] where scientists perform tests for whether colored jelly beans cause acne using a benchmark of *5%*. They test *20* different jelly bean colors and none are statistically significant – except the green one. The newspaper headline reads, "Green Jelly Beans Linked to Acne!" Since *20* tests are performed, on average *1/20 = 5%* false positives are expected. Accepting the "green jelly beans linked to acne" conclusion is a result of the multiple testing problem – if you run enough tests even when there is nothing, you will eventually "find" something by chance. There are a variety of procedures to adjust the benchmark to account for this. We will use the Bonferroni correction[36] because it is the simplest, it is widely used, and using more sophisticated methods would not change our

---

[34] "Statistical significance" and p-values have been the subject of recent criticism. See Ronald L. Wasserstein & Nicole A. Lazar, 'The ASA Statement on p-Values: Context, Process, and Purpose,' *The American Statistician*, 70.2 (2016), pp. 129–33. Nevertheless, the concept has served science well for a century in data problems very much like the one at hand.  Therefore, we intend to use p-values in that same tradition, while heeding the critics' warning: use p-values properly and with proper interpretation.  GB also use this technical phrase, as will be mentioned below.

[35] XKCD, 'Significant,' https://xkcd.com/882/. Accessed on January 12, 2024.

[36] See R. G. Miller Jr., *Simultaneous Statistical Inference* (New York: Springer-Verlag, 1991).

conclusions. The Bonferroni adjustment is applied by dividing the benchmark by the number of tests. In our case, we will have eighteen tests, therefore we will use the Bonferroni-adjusted benchmark of *0.05/18 = 0.0028*.[37]

### 2.2 Design of Tests and Hypothesis

The eighteen chi-squared goodness-of-fit tests performed with our data are as follows. We had five different test distributions (observed samples) of popularity/frequency data for our analysis: (1) GA, which is the focus of our study, (2) the works of Josephus, which are widely regarded as retaining historical reportage,[38] (3) historical novels, (4) the fictitious occurrences from Ilan-1 (Ilan-1f), and (5) GB's uniform distribution. We chose two different reference distributions: Ilan Vol. 1 Palestinian Jews (Ilan-1) and Ilan Vol. 3 Diaspora Jews (Ilan-3).

Our research hypothesis assumes that Ilan was correct to generally categorize the name occurrences from GA and Josephus as valid entries. Therefore, occurrences in GA and Josephus should "fit" Ilan-1 but not fit Ilan-3; additionally, the three hypothetical samples (historical novels, Ilan-1F, and the uniform distribution) would likely not fit into any historical reference distributions. Considering our aims and the data discussed previously, this resulted in predictions concerning eighteen tests. See Table 2 below.

| Reference Distribution | Variable | Gospels & Acts | Josephus | Novels | Ilan-1F | Uniform |
|---|---|---|---|---|---|---|

---

[37] In this paper, we are not using 0.0028 as a magic threshold upon which all p-values below it are statistically significant and all above are not. May it never be!  Indeed, as one reader correctly pointed out, we could select different numbers of tests to adjust this number. We merely offer it as an interpretive guideline on how to read the p-values shown below.  The specific number is not important. We could have done as few as four tests (GA popularity and origin vs. Ilan-1 and Ilan-3) and many more than eighteen. This could give Bonferroni-adjusted benchmarks in the range of 0.05/4 = 0.0125 to 0.05/100 = 0.0005.  The point is that "small" p-values indicate statistical significance and "large" do not and 0.0028 (or 0.0125 to 0.0005) will be a guideline and not a rule.

[38] Not that Josephus is immune from exaggeration and error. See Colin J. Hemer, *The Book of Acts in the Setting of Hellenistic History*, ed. C. H. Gempf (Winona Lake, IN: Eisenbrauns, 1990), pp. 97–98; Tal Ilan and Jonathan J. Price, 'Seven Onomastic Problems in Josephus' 'Bellum Judaicum',' *The Jewish Quarterly Review* 84.2/3 (1993), pp. 189–208.

| | | | | | | |
|---|---|---|---|---|---|---|
| Ilan-1 Palestinian Male Jews | Name Frequency | Fit | Fit | Not fit | Not fit | Not fit |
| | Name Origin | Fit | Fit | NA[39] | Not fit | Not fit |
| Ilan-3 Diaspora Male Jews | Name Frequency | Not fit | Not fit | Not fit | Not fit | Not fit |
| | Name Origin | Not fit | Not fit | NA | Not fit | Not fit |

*Table 2. The eighteen designed chi-square goodness-of-fit tests and our hypotheses.*

## 3. Results

### 3.1 Gospels & Acts, and Josephus

The p-values of the eighteen chi-square goodness-of-fit tests are shown in Table 3. Their output tables, along with some discussion of the technical details may be found in the Supplementary Materials.

| Ref. Distr. | Variable | Gospels & Acts | Josephus | Novels | Ilan-1F | Uniform (ACT) |
|---|---|---|---|---|---|---|

---

[39] We did not perform an analysis on name origin statistics on the historical novels, because the origin data for many of these names was not available in Ilan-1.

| | | | | | | |
|---|---|---|---|---|---|---|
| Ilan-1 | Name Freq. | 0.8556[40] | 0.0655[41] | $1.43 \times 10^{-14}$ | $2.20 \times 10^{-16}$ | $1.69 \times 10^{-15}$ |
| | Name Origin | 0.0034 | $1.09 \times 10^{-12}$ | NA | $2.20 \times 10^{-16}$ | $2.89 \times 10^{-6}$ |
| Ilan-3 | Name Freq. | $4.29 \times 10^{-13}$ | $2.20 \times 10^{-16}$ | $2.20 \times 10^{-16}$ | $2.20 \times 10^{-16}$ | $4.81 \times 10^{-5}$ |
| | Name Origin | $2.20 \times 10^{-16}$ | $2.20 \times 10^{-16}$ | NA | $2.20 \times 10^{-16}$ | 0.2819 |

*Table 3. Chi-square goodness-of-fit test results. Cell entries are p-values. The scientific notation, $1.19 \times 10^{-8}$, means the leading number (1.19 in this case) is divided by 10 to the 8th power. This means the decimal moves 8 places to the left, giving 0.0000000119 probability. The p-value $2.20 \times 10^{-16}$ is the smallest machine p-value without using extra precision. It is effectively zero. In the Uniform column, ACT stands for the GB "anonymous community transmission" hypothesis.*

Focusing on Table 3, we see that sixteen out of the eighteen tests matched our hypotheses from Table 2. Four of the p-values are greater than the Bonferroni-adjusted benchmark of *0.0028*, and the p-value for name frequencies in GA versus Ilan-1 is considerably high, providing evidence that GA fits Ilan-1 well. The two unexpected results are Josephus' name origin frequencies vs. Ilan-1 with p-value = $1.09 \times 10^{-12} < 0.0028$, which falls below our benchmark, and the name origin frequencies of ACT vs. Ilan-3 with a p-value *0.2819 > 0.0028*, which exceeds the benchmark.[42]

---

[40] Following GB, and in response to feedback, this test was GA vs. Ilan-1 with GA removed. Since Ilan-1 also includes Josephus, the test of Josephus vs. Ilan-1 also removed its name occurrences. This was also followed for the tests of GA and Josephus origin vs. Ilan-1. For all other tests against Ilan-1 and Ilan-3 the test distribution was not in Ilan-1 or Ilan-3 and so this adjustment was not made. For this test, with the name occurrences not removed from Ilan-1, the p-value was 0.9217.

[41] With the names not removed from Ilan-1, the p-value was 0.2143. There were 274, or 12.5% of the names removed from Ilan-1 for the test.

[42] The statistical explanation is as follows: For Ilan-1, there are 457 unique names among the 2185 – meaning there are a lot of replicates (around 5 occurrences per name). By contrast, for Ilan-3, there are 575 unique names among the 1227 – meaning there are far fewer replicates (around 2 occurrences per name). This means that sampling

In order to visualize the results of the hypothesis tests in Table 3, we have provided some figures which plot the reference distribution (Ilan) as a gray bar and the test distributions of GA, Josephus, the historical novels, Ilan-1F, and the Uniform distribution as grayscale 95% confidence intervals over the bar.[43] Figure 2 displays the popularity percentages and Figures 3 and 4 display the origin percentages. The way to interpret a single confidence interval is as follows: (i) the center of the interval is the actual proportion from the test distribution; (ii) the half-line above and below the center is the margin of error. If the interval goes through the top of the Ilan bar, then the test distribution is viewed as approximately fitting the reference for that single category. If the interval does not go through the top of the Ilan bar, then the test distribution is viewed as not fitting the reference for that single category. How close the interval is to the top of the bar matters; they can be "close" and yet "far off." Keep in mind that the confidence intervals are for visualization, which is subjective. They are calculated independently of one another, and actual fit is determined by formal testing, which accounts for all categories and closeness of fit simultaneously.

---

from the 575 (uniform) is "not too far" off from simply sampling from Ilan-3 itself, whereas sampling from the Ilan-1's 457 (uniform) is quite a bit different from sampling from Ilan-1's 2185.

[43] All confidence intervals are of the Wald type for proportions except the Uniform distribution, which had to be calculated using bootstrapping. See Supplementary Materials for details.
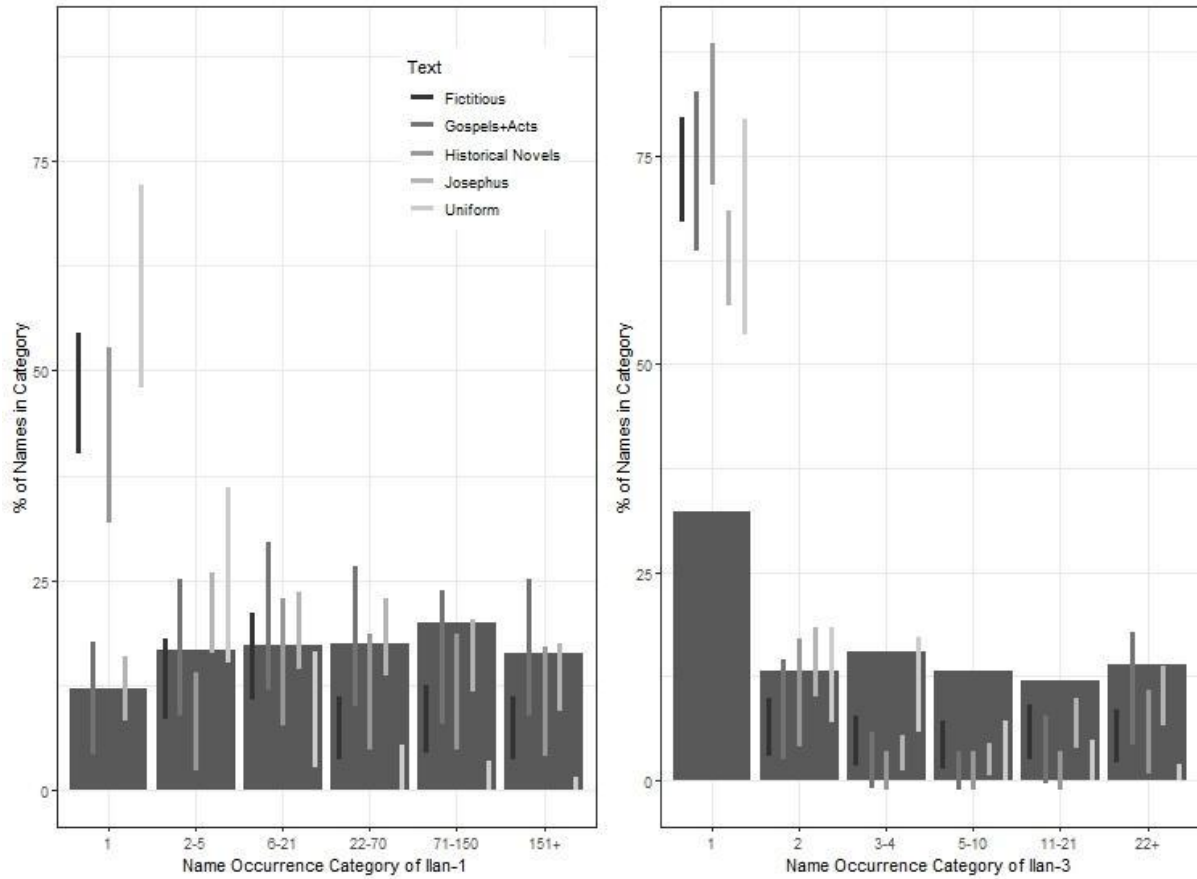
Figure 2. Popularity Statistics Comparing All Test Distributions with Ilan-1 and Ilan-3. The bar is the reference distribution. The bars are 95% confidence intervals. The match of Gospels & Acts, and Josephus, with Ilan-1 is seen with every confidence interval.
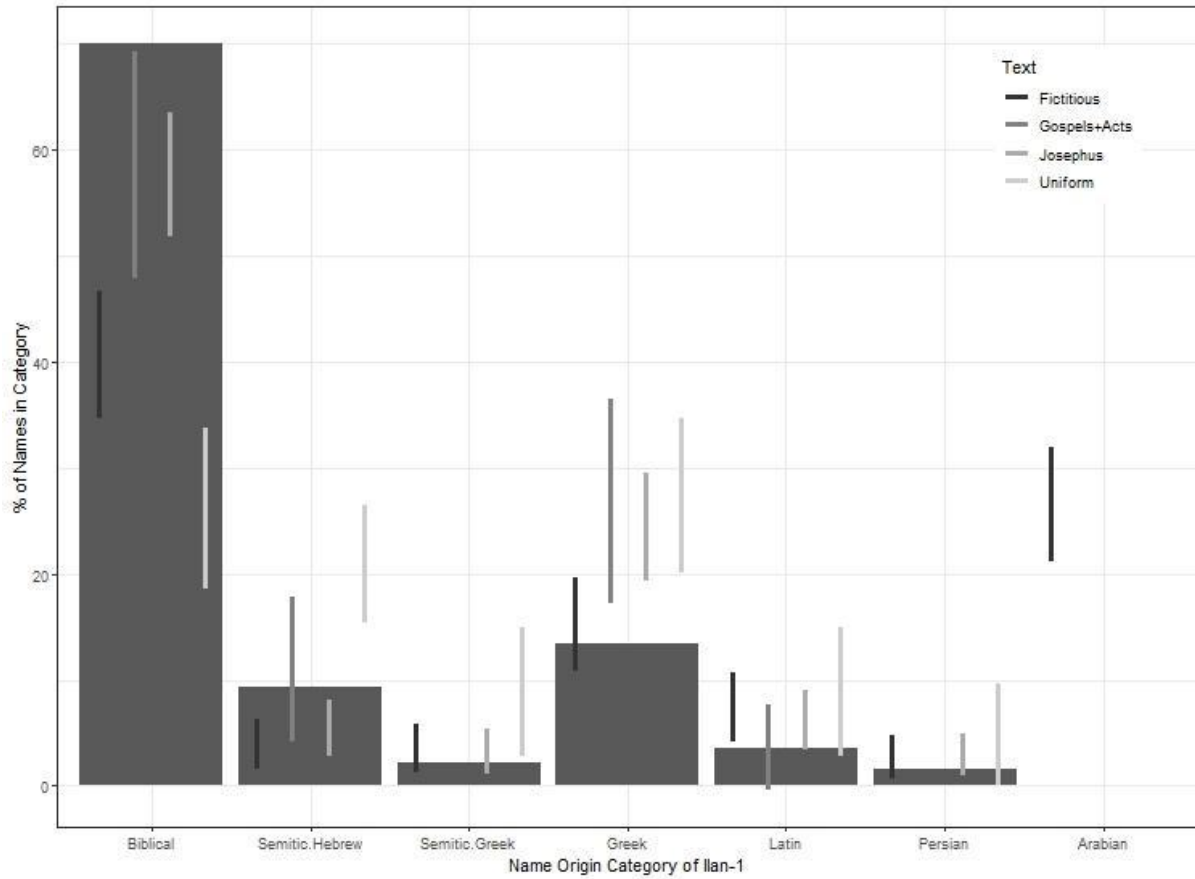
*Figure 3. Origin Statistics Comparing Gospels & Acts, Josephus, Ilan-1F, and Uniform with Ilan-1. The bar is the reference distribution of Ilan Vol. 1 Palestinian Males from 4 BCE to 73 CE. The bars are 95% confidence intervals.*
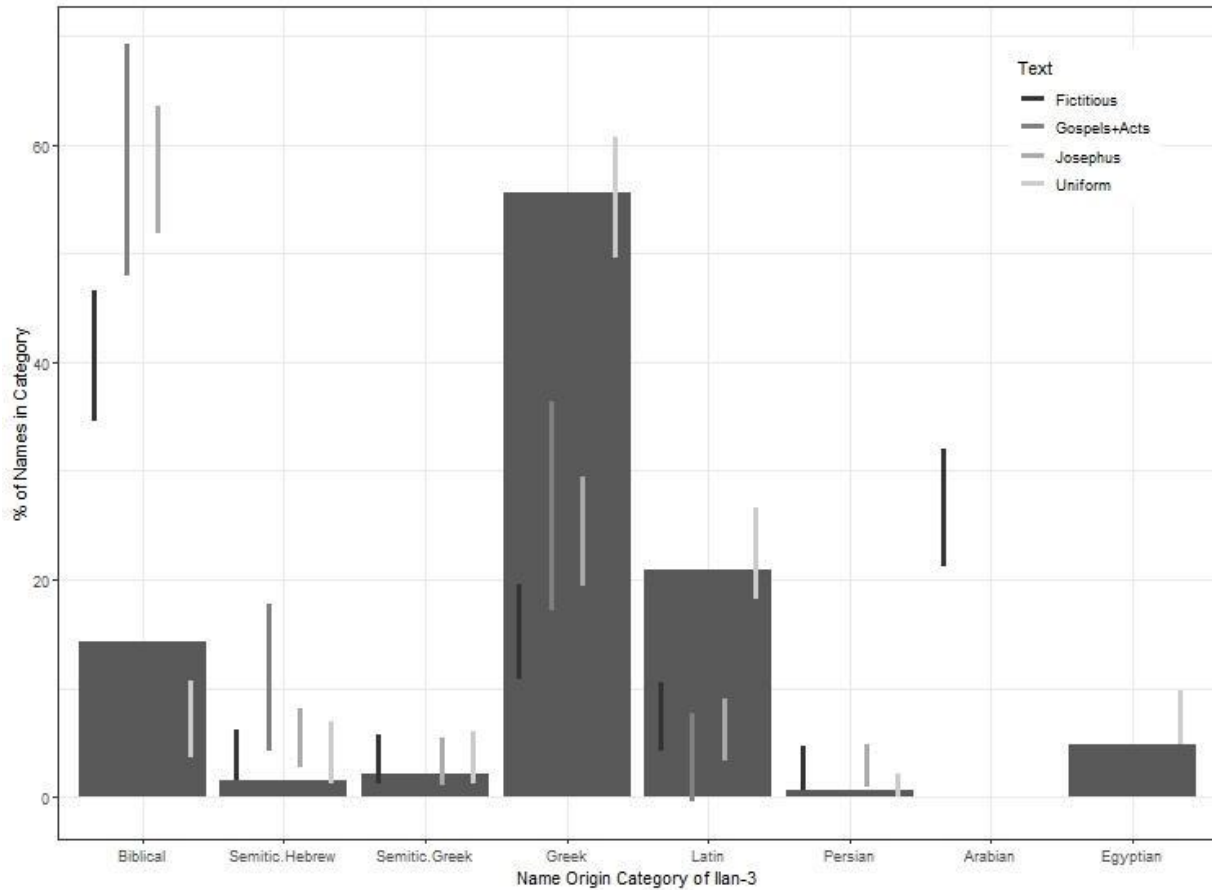
*Figure 4. Origin Statistics Comparing Gospels & Acts, Josephus, Ilan-1F, and Uniform with Ilan-3. The bar is the reference distribution of Ilan Vol. 3 Diaspora Males from 4 BCE to 73 CE. The bars are 95% confidence intervals. The mismatch between Gospels & Acts, and Josephus, with Ilan-3 is visible across multiple categories.*

Let us examine Figures 2, 3, and 4 and remark on their connection with the hypothesis tests from Table 3. In Figure 2, left side, we observe that all six intervals of GA and Josephus overlap Ilan-1, consistent with the high p-values in Table 3. GA is slightly low on occurrences of 71-150 times in Ilan-1, likely due to the fewer occurrences of Eleazar, while Josephus is high in occurrences of 2-5 and low in occurrences of 151+. Turning to the right side of Figure 2, only two of the six GA and Josephus intervals overlap Ilan-3, reinforcing the extremely low p-values in Table 3.

Moving to the origin statistics, we start with Figure 4 vs. Ilan-3. Again, neither GA nor Josephus visibly "fit." This is most prominent in that the biblical percentage of names is way too high whereas the Greek percentage is way too low. Semitic Hebrew names are also too high,

balanced with Latin too low. Lastly, neither GA nor Josephus have any Egyptian names and therefore their percentage is 0 with no confidence interval able to be calculated.  This is a clear case of no fit, as corroborated with the GA and Josephus p-values of $2.20 \times 10^{-16}$ in Table 3.

However, in Figure 3 things are not so clear. As previously, GA and Josephus track similarly. However, GA is slightly low on Biblical and high on Greek, but not much.  Josephus is a little lower on biblical and higher on Greek and Latin, and also a tad below on Semitic Hebrew. Notice also that the GA confidence intervals are always wider than Josephus. This is because Josephus' sample size is 274 whereas GA is only 82, providing less confidence and therefore wider intervals. All these factors combine to produce a figure that does not clearly show the difference in the magnitude of fit to the naked eye. But the p-value of the GA is 0.0034 whereas the Josephus p-value is $1.09 \times 10^{-12}$ in Table 3. This illustrates the reason why statisticians do not rely on subjective interpretations of graphs to draw their conclusions, but rather draw conclusions from appropriately chosen and applied statistical tests. One surprising result of our analysis, therefore, is that Josephus' name origin data statistically does not fit Ilan-1 due to a proportionately high occurrence of Greek and Latin names versus biblical and Semitic Hebrew names.

Several factors likely contribute to this. Ilan's database relies on epigraphic as well as literary sources, including ossuaries from the region of Jerusalem, which show a preference for biblical over Greek names. Conversely, Josephus' audience, focus, and literary milieu results in a greater preference for Greek over biblical or Semitic names of Jewish persons. This is exemplified in his generally more Hellenized orthography and formal Greek case suffixes, while GA show a preference toward less official orthography coinciding, as Ilan notes, to common pronunciation.[44] Discussing the general prominence of Greek names in literary sources, Ilan observes, "The literary texts record the lives of the more affluent and worldly, who chose names for their children from the culture by which they were influenced. The epigraphic record documents a wider variety of people and shows that this apparent trend in the literary sources was much less popular in the wider population."[45] Compared to Josephus, GA provides a slightly

---

[44] Ilan, *Part I*, p. 18.

[45] Ilan, *Part I*, p. 40.

more unvarnished, on-the-ground situation that more accurately reflects general naming practices, and the fitness of GA as opposed to Josephus in this instance is not merely due to the smaller sample size.[46]

### 3.2 Hypothetical Test Distributions

Turning again to Figure 2, the occurrences from the historical novels and Ilan-1F vs. Ilan-1 generally do not fit, especially in the over-representation of infrequent names. This data runs counter to the argument from GB and is in accordance with observations made elsewhere that name rarity is typically over-represented by inventors of names.[47] Although we combined *Ben Hur* and *The Spear* into a single sample for our analysis, it is worth noting that their p-values fall significantly below our threshold when considered separately as well. *Ben Hur* has a p-value of $9.071 \times 10^{-6}$ and *The Spear* a p-value of $2.831 \times 10^{-15}$.

As for the uniform distribution vs. Ilan-1, there is an over-representation in frequency categories "1" and "2-5", with severe under-representation in "71-150" and "151+". This is not surprising since each name was represented only once. As the p-values in Table 3 demonstrate, testing GA and these hypothetical samples against Ilan-1 demonstrates that at least for a male Jewish Palestinian population (4 BCE–73 CE), the number of name occurrences in GA has been shown to fit the reference distribution, while these other purported test distributions failed to fit.

---

[46] The literary setting also influences the GA sample (e.g., "Peter" instead of "Cephas"), which is also tilted in favor of having higher Greek occurrences than expected due to the historical situation described in Acts 6. Note that adding the six deacons from Acts 6:5, as GB also observe ('Name Popularity,' p. 185), significantly impacts the ratio of GA's Greek to biblical names in relation to the Ilan-1 reference distribution. Without the six deacons from Acts 6:5 (Stephen, Philip, Procorus, Nicanor, Timon, and Parmenas), the p-value of GA's origin statistics versus Ilan-1 is .0587, exceeding the standard threshold for statistical significance even apart from the Bonferroni correction. The historical situation described in Acts readily supplies the reason why these Jewish males would have disproportionately borne Greek versus biblical names: they were specifically chosen to represent Hellenistic Jews (Acts 6:1-5). As an aside, combining Semitic-Hebrew and Semitic-Greek into a single category results in GA p-values of .065 (deacons excluded) and .0037 (deacons included).

[47] See Van de Weghe, 'Name Recall,' pp. 100–01.

Let us now briefly consider the 26 uncontested GA name occurrences as cataloged by GB.[48] When considering these against a reference distribution of Ilan-1, they have a p-value of 0.1369. Despite the smallness of this sample size, this is lower than the p-value of the contested names alone (p-value = .7450) as well as the p-value of the complete GA sample: p-value = 0.8543). Adding the contested names so as to consider the full sample of GA names results in at least as good of a historically situated fit. This would be expected if the contested name occurrences are authentic.[49]

Even if, contrary to the evidence we possess,[50] accurate and intricate naming patterns were achieved by a fictionalizing author of antiquity, there are several additional reasons that the authors of GA were not likely to have done so. First, in the apocryphal material any hint of Jesus' Palestinian environment is overshadowed by mythological concerns and persons while knowledge of Palestinian persons becomes confused (e.g., thinking Peter and Cephas are two different people).[51] Second, an unusual name like Zenon, a more popular Western Diaspora name like Samuel, or a New Testament name like Zacchaeus is occasionally added, but the general trend reflects a disinterest in capturing any kind of realistic onomastic snapshot.[52] Third, to assume that a later fictional writer would diligently examine a source like Josephus or other material from Palestine, take a rough inventory of name popularity, and attempt to reproduce it assumes, against available evidence, a hyper-sensitive focus on names that no critical NT reader considered meaningful for two millennia until Richard Bauckham began the research that resulted in *Jesus and the Eyewitnesses*. These three observations align with the published data

---

[48] Gregor and Blais, 'Name Popularity,' pp. 185–86.

[49] We thank Willem Jan Blom for bringing this point to our attention.

[50] Van de Weghe, 'Name Recall,' pp. 101–04.

[51] Peter and Cephas are regarded as different persons in *The Epistle of the Apostles* (see Simon Gathercole, *The Apocryphal Gospels* (Westminster, London: Penguin Books, 2021), p. 261. A mythological focus is evident, for example, in *The Sophia of Jesus Christ*, *The Gospel of Judas*, *The Coptic Gospel of the Egyptians*, and *The Gospel of Mary.*

[52] For a survey of named characters in apocryphal gospels, see the lists of *dramatis personae* in Gathercole, *The Apocryphal Gospels*, pp. 4, 31, 48, 147-48, 155, 163, 194, 207-11, 254, 261, 304, 361.

on personal name retention, which suggests that names, being arbitrary and easily forgettable, are among the least integrated and lasting pieces of information in the recollection of stories.[53]

In this regard, we must also consider the unique historical setting in which our data sample occurs. While naming frequencies in the Greco-Roman world were generally distributed widely and evenly, the uniquely high occurrences of only several Hasmonean names make it much harder for a fictional author to achieve a level of apparent authenticity.[54] This explains why, surprisingly, the name origin frequencies of "anonymous community transmission" (uniform) vs. Ilan-3 show a fit with a p-value of 0.2819 with a sample size of 52 (GA contested name only size) and 0.2026 with a sample size of 82 (full GA size), which indisputably exceeds the benchmark.[55] This demonstrates that while it is possible that an "anonymous community transmission" (uniform) distribution may fit some real reference distributions (Ilan-3), it does not fit others (Ilan-1).  This provides evidence that the fit signal from GA to Ilan-1 can be statistically distinguished from the lack of fit signal of GA to Ilan-3.

This situation in Ilan-1, wherein a relatively high number of persons bear relatively fewer names, creates further patterns of authenticity in GA that are more difficult to quantify. Note, for example, that from the 28 cases in which a GA author disambiguates a person's primary name (e.g., Simon *the Leper*, Judas *Iscariot*, etc.), 75% of these belong to persons bearing one of

---

[53] See the discussion in Van de Weghe, 'Name Recall,' pp. 95-96.

[54] As Ilan observes, 20.7% of the name pool served 73.4% of the male population (*Part I*, p. 5). To illustrate the uniqueness of this, consider the top two names in Ilan-1. These belong to more than 16% of all males. The top two Greek names in Coastal Asia Minor belong to under 5% of the male population (1,849 out of 39,477 occurrences; P. M. Fraser and E. Matthews, eds., *The Lexicon of Greek Personal Names, Volume V.B, Coastal Asia Minor: Caria to Cilicia* (Oxford: Clarendon, 2013), pp. xxx-xxxi.

[55] Part of the reason is due to the uniformity of name occurrences within the Ilan-3 population compared to Ilan-1. For Ilan-1, there are 457 unique names among 2185 occurrences, meaning there are more replicates.  This averages to about 5 occurrences per name, but in reality some occur with high frequency while many are rare (only one occurrence). By contrast, for Ilan-3 there are 575 unique names among 1227 occurrences, meaning there is far less high/low frequency variation. That fact alone does not result in a fit with uniform distribution frequencies (p-value = 2.20×10-16), but that combined with the particular mix of frequency-origin combinations can.

the top five names and 89% to persons bearing one of the top 12 names.[56] This reflects, again, a "situation on the ground" wherein men bearing the most common names would naturally need to be distinguished from one another; it is contrary to the practice in De Wohl's *The Spear*, which shows a tendency toward qualifying names regardless of their popularity.[57] In GA, the situation becomes even more precise. For example, Lazarus "of Bethany" is disambiguated (John 11:1); this name occurs widely in the general population, but not in GA itself. The common name "Jesus" is consistently disambiguated throughout the Gospels in the public speech of characters but not by the Gospel authors themselves in segments of narration, as would be expected under an authentic information-retention scenario.[58] Much like the less official orthography of GA, the naming practices at a meticulous level favor a model that allows for natural name retention on a minute scale, whereas the data observed in our analysis does not support a model incorporating the achievements of even a well-researched forger of this information.

## 4. Remarks on Gregor and Blais

Having completed our analysis of the data, we now explicitly address the work of GB in the following remarks. These highlight why our findings were distinct from GB and where GB's statistical analysis falls short.

---

[56] This number excludes qualifiers due to titles (e.g., Caiaphas the high priest, king Herod, etc.); see our GA data. Including titles brings the total number of qualified name occurrences to 36. Even with this number, 75% of qualifiers are given to names ranking in the top 12 of Bauckham's table 6 in *Jesus and the Eyewitnesses*, p. 84.

[57] Of the 26 qualified names in *The Spear*, 50% of them have under ten occurrences in Ilan-1 (all besides Nathan and Hillel have five occurrences or less): Boz, Zadok, Aaron, Ephraim, Achim, Oziah, Amram, Nathan, Mordecai, Josaphat, Nicodemus, Baruch, and Hillel. When *The Spear* includes a qualifier with a popular first name, it is often because the entire name, including the qualifier, is taken from GA (e.g., Simon bar Jonah, Judah son of Alphaeus, Yohanan son of Zebedee, Annas the high priest, Joseph of Arimathea, Jesus bar Joseph, Yohanan the Baptist, Judas from Kerioth). We thank Lydia McGrew for drawing our attention to this trend.

[58] This is because the authors clearly knew which Jesus they were talking about, while in the general public his name would require a disambiguation. Peter Williams details this phenomenon in *Can We Trust the Gospels?* (Wheaton, IL: Crossway, 2018), pp. 71–75.

(4.1) The primary argument of GB rests on their Figures 1 and 2 which show 95% confidence intervals of name frequencies, somewhat like our own Figures 2, 3, and 4. The first major problem with this approach is that their conclusion is subjective because it rests on their interpretation of their graphs. We disagree with their interpretation and would actually draw the opposite conclusion from their graphs.[59]

(4.2) The second major problem with GB basing their argument on their Figures 1 and 2 is that they implicitly assert a statistical hypothesis test[60] with a null hypothesis of "anonymous community transmission", or random name selection.[61] The problem is that they went contrary to standard statistical practice and reversed their null and alternative hypotheses.[62]

---

[59] GB write, "Most observed Gospels-Acts numbers of name occurrences fit inside the confidence interval…" (p. 24). In fact, not just "most", but 25 out of the 26 names shown fit the confidence intervals. We contend that this is indicative of fit, not lack of fit. Reducing subjectivity is one of the reasons for hypothesis testing in the discipline of statistics, as demonstrated in Section 3.

[60] Statistical hypothesis tests consist of two hypotheses: the null hypothesis and the alternative hypothesis. They are structured such that the null hypothesis is the status quo and assumed to be true unless the preponderance of evidence from the data indicates otherwise, in which case the alternative hypothesis is concluded. In our paper, "the preponderance of evidence" is given by the benchmark of 0.05, divided by the 18 tests, which is 0.0028. This is an objective criterion: when a p-value is above 0.0028 the data supports the null hypothesis, when the p-value falls below 0.0028 the null hypothesis should be rejected in favor of the alternative, in the absence of countervailing evidence.

[61] GB never explicitly state a formal statistical hypothesis test. However, some (not all) of their language, approach, and Supplementary Materials refer to an implied hypothesis test. This is most clearly seen with their use of the technical phrase "statistically significantly different" on p. 197 of 'Name Popularity,' which refers to rejecting a null hypothesis in favor of the alternative hypothesis. "Any combination of at least some of the contested characters being historical and the invention of fictitious characters' names with at least some information about name popularity would result in a name popularity distribution corresponding to the contemporary population distribution more closely than the distribution generated on the most extreme scenario. Because the sample size of Gospels-Acts name occurrences is so small that the observed Gospels-Acts name popularity distribution is already not *statistically significantly different* from [anonymous community transmission]" (emphasis ours).

[62] In statistical goodness-of-fit testing, of which the chi-square goodness-of-fit test is the appropriate test for our data type, the null hypothesis is 'the data fit the reference distribution.' In our case, the reference distribution must be the historical Ilan, or something of that sort. This is simply how the tests work, and GB have silently

(4.3) Third, using the method of GB, a sample of size 53 (or 52) from many other text's distributions would draw the same conclusion of 'fit' to the uniform distribution, i.e. "anonymous community transmission." To see this, look at their Figure 1: the top of the gray bar passes through, or at least touches, all of the confidence intervals shown. Because of this, GB argues that GA fits the uniform distribution. By contrast, in their Figure 2, not all of the Josephus confidence intervals go through the gray bar, therefore Josephus does not fit the uniform distribution, they argue. The problem is that if a sample of this size were selected from Josephus, then the confidence intervals of Figure 2 would correspondingly widen and look like Figure 1, switching the conclusion of "no fit" to "fit." No one thinks Josephus' writings were "not statistically significantly different"[63] from anonymous community transmission – yet GB's methodology would conclude this if he had written less! This is a flaw in their methodology.[64] This fallacy arises from switching the null and alternative hypotheses as in the above remark.

(4.4) Fourth, the confidence interval method employed by GB in Figures 1 and 2, as well as in our own Figures 2, 3, and 4, assume independence between names, which is false. Such confidence intervals are used because they are easy to calculate, relatively easy to explain, and make a great picture to illustrate the point. Nevertheless, the proper use of the figures is for illustration only. In reality, the proportion of one name is dependent upon the proportions of all

---

switched their null and alternative hypotheses. The reason the null must be the reference distribution is that it is well-defined, whereas the possible alternative distributions are uncountable. They would certainly include "anonymous community transmission," but which one?  Defining this "anonymous community transmission" distribution is problematic.  See Supplementary Materials for some discussion. GB opt for a uniform distribution based on the 451 unique names in their adjusted Ilan-1 plus the 3 they returned to Ilan-1. Why not use a mixture uniform distribution with two classes from that set? Why not use the set of just contested/uncontested New Testament names? The point is that there are an uncountable number of alternative possibilities and that is precisely the reason that goodness-of-fit testing places the reference distribution in the null hypothesis.

[63] Gregor and Blais, 'Name Popularity,' p. 197.

[64] Additionally, Ilan-1F demonstrates that forgers often incorporated names into their narratives that occur nowhere in Ilan-1; their consideration of only Palestinian names is not justified by the small sample of occurrences from the *Protoevangelium of James*. Regardless, the p-value of 52 uniform ACT occurrences vs. Ilan-1 is significantly lower ($1.69 \times 10^{-15}$) than the p-value of the 52 contested GA occurrences vs. Ilan-1 (0.7450).

of the other names, which is not accounted for in the confidence interval calculations. However, this dependence relationship is precisely accounted for in the chi-squared goodness-of-fit test we have used earlier in this paper.

(4.5) Fifth, GB assert that names occurring only once in GA provide no information about whether GA fits the historical distribution.[65] However, as we have shown above, grouping individuals in order to meet the test conditions is a common issue in goodness-of-fit testing, and the information for every name can be accounted for by binning them into appropriate groups, as we have done.

(4.6) Sixth, GB's Section 5 directly addresses the issue of rare names. This is a welcome step. Rare names are those which occur only once in Ilan-1.[66] We overlook the fact that GB previously claimed that names occurring only once in GA were "without any information about name popularity."[67] However, the method used in this section has substantial problems. To start, they reversed the null and alternative hypotheses which they used in Section 4. Although we agree that this is now the correct configuration of the hypotheses, it is an outright inconsistency.[68] The main issue is that GB's conclusion in Section 5 is highly sensitive to the number of rare names in the list, and we dispute their list, which stands at four (Aeneas, Agabus, Bartimaeus, and Timaeus). In their handling of the data, GB removed 26 "attested"[69]

---

[65] "Therefore, most of the Gospels-Acts data is entirely uninformative and cannot support Bauckham's thesis. It is statistically indistinguishable from 'white noise' generated entirely randomly without any information about name popularity," p. 191.

[66] In their article, GB define rare names as "attested only once among first-century Palestinian Jews" ('Name Popularity,' pp. 174, 198). If any readers wish to use GB's Ilan file, as we have done, note that it defines rare as either attested once (262 names), or twice (79 names, except "Yeshab" appears to have inadvertently not been marked as rare).

[67] See the fifth remark in this section.

[68] See our remarks in 4.2 and 4.3. In GB's Section 4 their null hypothesis was anonymous community transmission and the alternative was Ilan-1. Now, in GB's Section 5 their null hypothesis is Ilan-1 and the alternative is anonymous community transmission.

[69] Removing the 26 "attested" names from the GA sample is built on questionable assumptions that we do not agree with. Foremost, it alters the GA dataset to favor their conclusion (in this case, removing names from a

name occurrences (21 names) from Bauckham's original 79 (45 names) to obtain their 53 (32 names) unattested occurrences.[70] The names removed include two rare names (Qaifa and Toma). According to their calculations, the "probability of getting [four or fewer] names … is only 1.1%!"[71] According to GB's argument, the table below shows the number of rare names and the probability of that many or fewer occurring, if the Ilan-1 distribution was true.[72] It can be seen that their conclusion is extremely sensitive to the number of rare names. GB's list gives a 1.1% probability, which would jump to 7.3% if just the two names were not excluded (Table 4). Moreover, if the Palestinian Acts deacons were added, as we believe they should be, it includes four more "rare" names and would bring the list to eight or ten rare names at 24% or 50% of occurring (Table 4), depending on whether we also add the two contested rare names, Qaifa and Toma.[73]

---

dataset and then concluding that the dataset is too small to determine statistical significance). Further, removing "attested" names is justified by a narrow focus on several overstatements made by Bauckham (Gregor and Blais, 'Name Popularity,' pp. 179, 180, 181; see also, Van de Weghe, *Living Footnotes*, p. 34), while a broader reading of Bauckham could allow all the GA data to be considered as part of a cumulative argument for authenticity. Further, the notion of "attested" is based on ad hoc considerations; GB do not try to argue that GA authors relied on any of these sources for their name statistics, for example, but merely that we can know from other sources that these names were historical persons. But the subjectiveness of whether we take such external attestations seriously is epitomized in the implicit assumption that names in the fragments of Papias' writings should be taken at face value while the names in GA should be brought into question. Lastly, this creates a very uncomfortable scenario in which historically-attested GA names function to *undermine* the historicity of names in GA as part of GB's broader argument.

[70] The reason the 45 names only reduces to 32 names when 21 names are removed is that 8 names have separate occurrences in both lists.

[71] Gregor and Blais, 'Name Popularity,' p. 198.

[72] See GB Figure 3. The actual probabilities are not shown in the article, but they are closely approximated by a binomial distribution with n=53 and p=0.53, which are the results shown, including the 1.1%. See Supplementary Materials.

[73] Since we remove Bartimaeus, our sample includes nine rare GA names. Reflecting Ilan's orthography, they are: Hagaba, Timaeus, Qaifa, Parmenas, Procharus, Stephanus, Toma, Timon, and Aeneas. Aeneas is attested once in Ilan-1 outside of GA, so including this in the rare category depends on how one categorizes the GA occurrence.

| # rare | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|--------|-----|-----|-----|-----|----|----|----|----|
| % | 0.3 | 1.1 | 3.2 | 7.3 | 14 | 24 | 36 | 50 |

*Table 4. Probability of obtaining numbers of rare names. This table is based off GB's Figure 3. For example, # rare = 6 and %=7.3 means that there is a 7.3% probability that 6 or fewer rare names would be obtained in a sample of size 53 selected randomly from GB's pool of 2,582 occurrences.*

## 5. Conclusion

In conclusion, GB's analysis suffers from a lack of proper statistical methodology, thereby introducing critical errors into their analysis. In this paper, we have corrected their missteps and have attempted to analyze as much data as possible with the best available statistical methods. We have shown how chi-square goodness-of-fit testing establishes that naming patterns in the Gospels and Acts fit into their historical context well, and that they fit statistically significantly better than fictitious works from antiquity and well-researched historical novels from the modern era. We have considered the GA name sample in light of two reference distributions and also discussed name origin data, sample size, and less quantifiable observations about name disambiguation and naming practices within invented materials. The evidence suggests that GA accurately retained personal names – those unmemorable pieces of personal information – at a remarkably high level. In the words of Richard Bauckham, "all the evidence indicates the general authenticity of the personal names in the Gospels."[74]

---

[74] Bauckham, *Jesus and the Eyewitnesses*, p. 84.