

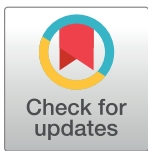
RESEARCH ARTICLE

Published estimates of group differences in multisensory integration are inflated

John F. Magnotti*, Michael S. Beauchamp

Department of Neurosurgery and Core for Advanced MRI, Baylor College of Medicine, Houston, Texas, United States of America

* magnotti@bcm.edu



Abstract

A common measure of multisensory integration is the McGurk effect, an illusion in which incongruent auditory and visual speech are integrated to produce an entirely different percept. Published studies report that participants who differ in age, gender, culture, native language, or traits related to neurological or psychiatric disorders also differ in their susceptibility to the McGurk effect. These group-level differences are used as evidence for fundamental alterations in sensory processing between populations. Using empirical data and statistical simulations tested under a range of conditions, we show that published estimates of group differences in the McGurk effect are inflated when only statistically significant ($p < 0.05$) results are published. With a sample size typical of published studies, a group difference of 10% would be reported as 31%. As a consequence of this inflation, follow-up studies often fail to replicate published reports of large between-group differences. Inaccurate estimates of effect sizes and replication failures are especially problematic in studies of clinical populations involving expensive and time-consuming interventions, such as training paradigms to improve sensory processing. Reducing effect size inflation and increasing replicability requires increasing the number of participants by an order of magnitude compared with current practice.

OPEN ACCESS

Citation: Magnotti JF, Beauchamp MS (2018) Published estimates of group differences in multisensory integration are inflated. PLoS ONE 13(9): e0202908. <https://doi.org/10.1371/journal.pone.0202908>

Editor: Mark W. Greenlee, Universitat Regensburg, GERMANY

Received: January 22, 2018

Accepted: August 10, 2018

Published: September 19, 2018

Copyright: © 2018 Magnotti, Beauchamp. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: All reported data and source code for reported simulations are available online: <http://www.openwetware.org/wiki/Beauchamp:DataSharing>.

Funding: This research was supported by NIH R01NS065395 to MSB. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing interests: The authors have declared that no competing interests exist.

Introduction

Since different sensory modalities carry distinct sources of information about the world, integrating them provides a more reliable picture of the world. The underlying computations and behavioral manifestations of this multisensory integration have received increasing attention. Researchers have developed a variety of measures to assess multisensory integration, focused on a change in perception due to the addition of a second sensory modality. For instance, in speech perception, integrating auditory information from the talker's voice and visual information from the talker's face enhances speech recognition accuracy [1–3]. Perhaps the most common measure of multisensory integration is the McGurk effect, an illusion in which incongruent auditory and visual speech are integrated to produce a new percept that matches neither of the component modalities. The original report of the McGurk effect [4] has been cited thousands of times (Fig 1A) and the illusion is used as a textbook example of

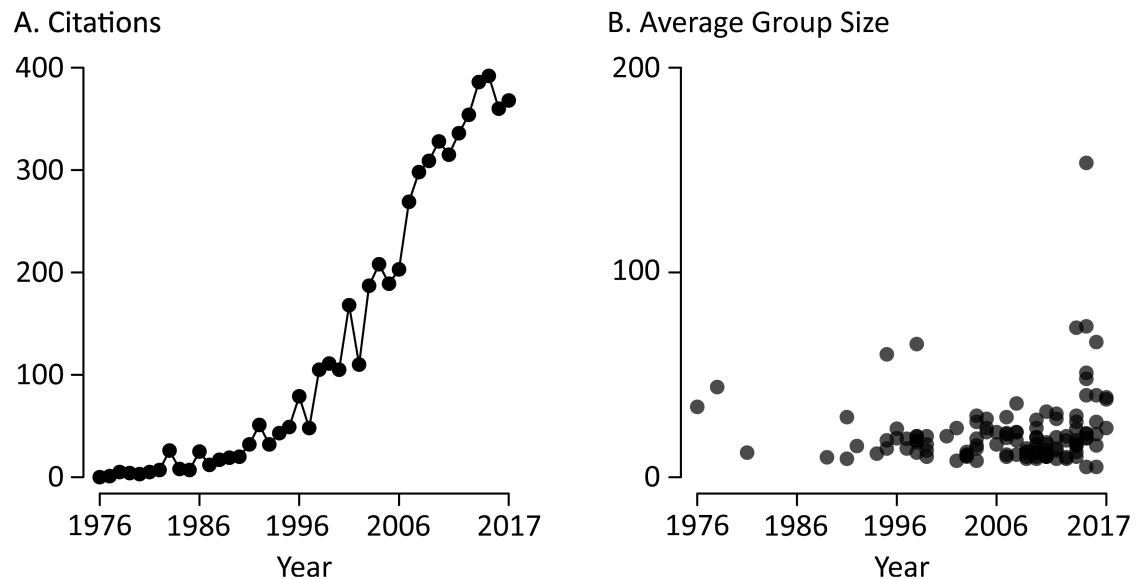


Fig 1. The McGurk effect in the scientific literature. (A) The number of citations per year of the McGurk and MacDonald Nature paper describing the illusion, since initial publication in 1976. (B) An analysis of the sample size (defined as the number of participants in an experimental group) in published papers on the McGurk effect (119 total papers, 262 total experimental groups).

<https://doi.org/10.1371/journal.pone.0202908.g001>

multisensory integration. As a behavioral assay, susceptibility to the McGurk effect has been used to argue for important differences in multisensory integration between genders [5], in atypical human development [4, 6], across the lifespan [7], in mental health disorders [8, 9], and between cultural/linguistic groups [10–12]. Amidst the enthusiasm for using the McGurk effect to study such between-group differences, recent work has highlighted large within-group variability in susceptibility to the illusion using relatively homogeneous subject pools [13, 14]. However, most studies of the McGurk effect test a relatively small number of participants, with only a single published study reporting a group size of greater than 100 [15]. We demonstrate that high within-group variability, combined with small sample sizes and the common scientific practice of publishing only significant results, leads to inflated estimates of group differences in multisensory integration. These results explain why follow-up studies have failed to replicate published reports of large differences in the McGurk effect between cultures [15], between genders [13], and between typically-developing controls and children with developmental disorders [16]. The proliferation of published reports with inaccurate estimates of group differences has led to stark conflicts in the literature, with different studies examining similar populations reporting completely opposite results. This situation makes it difficult to test scientific theories or develop therapies for patients with multisensory deficits [17]. Studies of multisensory integration must increase sample sizes by an order of magnitude to produce reliable, low-variance estimates of group differences.

Results

To demonstrate how small studies can produce inflated effect estimates in the presence of a publication filter, we modeled the consequences of studying group differences in multisensory integration with at a variety of sample sizes. This modeling effort is made possible by a recent large-scale, in-person behavioral study, allowing an accurate sample of the true variability within a large population of homogeneous, young healthy college undergraduates [13]. For

other situations, such as comparisons with clinical populations, the within-group variance is expected to be at least as large as for college undergraduates. Therefore, these simulations provide a lower bound on the sample sizes necessary to study group differences in susceptibility to the McGurk illusion. Using empirical data avoids assumptions, such as population normality, that are demonstrably false for the McGurk effect [13].

Effect inflation measured with a known true effect

To examine how sample size can influence experimental results, we used bootstrapped data from a large behavioral study and simulated population differences in McGurk susceptibility, defined as the mean percentage of fusion responses to McGurk stimuli across all possible stimuli and subjects (see *Methods* for details). In the first example, we created two populations: population A has susceptibility 45% while population B has susceptibility 55%, and thus a mean difference of 10% (Fig 2A). How likely is it for a given experiment to accurately estimate this group difference? Fig 2B shows an example experiment in which 150 subjects are sampled from each population for a total sample size of 300. With this large sample size, susceptibilities are estimated at $48\% \pm 3\%$ (standard error of the mean) for population A and $60\% \pm 3\%$ for population B, resulting in a difference estimate of $12\% \pm 9\%$ (95% confidence interval), a reasonable approximation to the true population difference of 10%.

However, a sample size of $N = 300$ is much larger than that used in published studies of multisensory integration. Fig 2C shows an example experiment with a more typical sample size of $N = 30$. This experiment results in estimated susceptibilities of $36\% \pm 9\%$ for population A and $67\% \pm 8\%$ for population B. The estimated population difference in this experiment is $31\% \pm 24\%$, greatly inflated from the actual difference of 10%.

This effect inflation occurs when experimenters only consider the results of experiments that result in a significant population difference, usually defined as a between-groups *t*-test producing $p < 0.05$. This criterion is deeply embedded in the scientific process, most often in the manuscript preparation and submission phase in which only significant results are included [18–20]. Fig 2D shows the result of implementing this significance filter on thousands of simulated experiments with a sample size of 300. Of these experiments, 67% result in significant differences between populations and are "published"; the other 33% are discarded. The discarded studies have population difference estimates that are always *smaller* than the true difference (mean of 5%; 100% of estimates less than 10%; gray bars in Fig 1D) while the significant studies have population difference estimates that are usually *greater* than the true difference (mean of 13%; 77% of estimates greater than 10%; blue bars in Fig 2E). A weighted average of the significant and discarded experiments recovers the true population difference of 10% ($13\% \times 0.67 + 5\% \times 0.32 = 10\%$), but examining only the significant experiments biases the population difference estimate upwards by a factor of 1.3.

Fig 2E shows the results of implementing a significance filter on thousands of simulated experiments with sample size of 30. Only 13% of experiments result in significant differences between populations; the remaining 87% are discarded. The discarded studies have population estimates that are close to the true difference (mean of 7%; gray bars in Fig 2E), while the included studies have population difference estimates that are much greater than the true difference (mean of 30%; blue bars in Fig 2E). Considering only significant experiments inflates the population difference estimates by a factor of 3.

A key point is that small sample sizes not only inflate the mean effect estimate (when only significant results are published), but also increase effect estimate variance from study to study (regardless of the publication filter). For any particular $N = 30$ study with a true population difference of 10%, the difference estimate will vary from -16% to 36% (two standard deviations

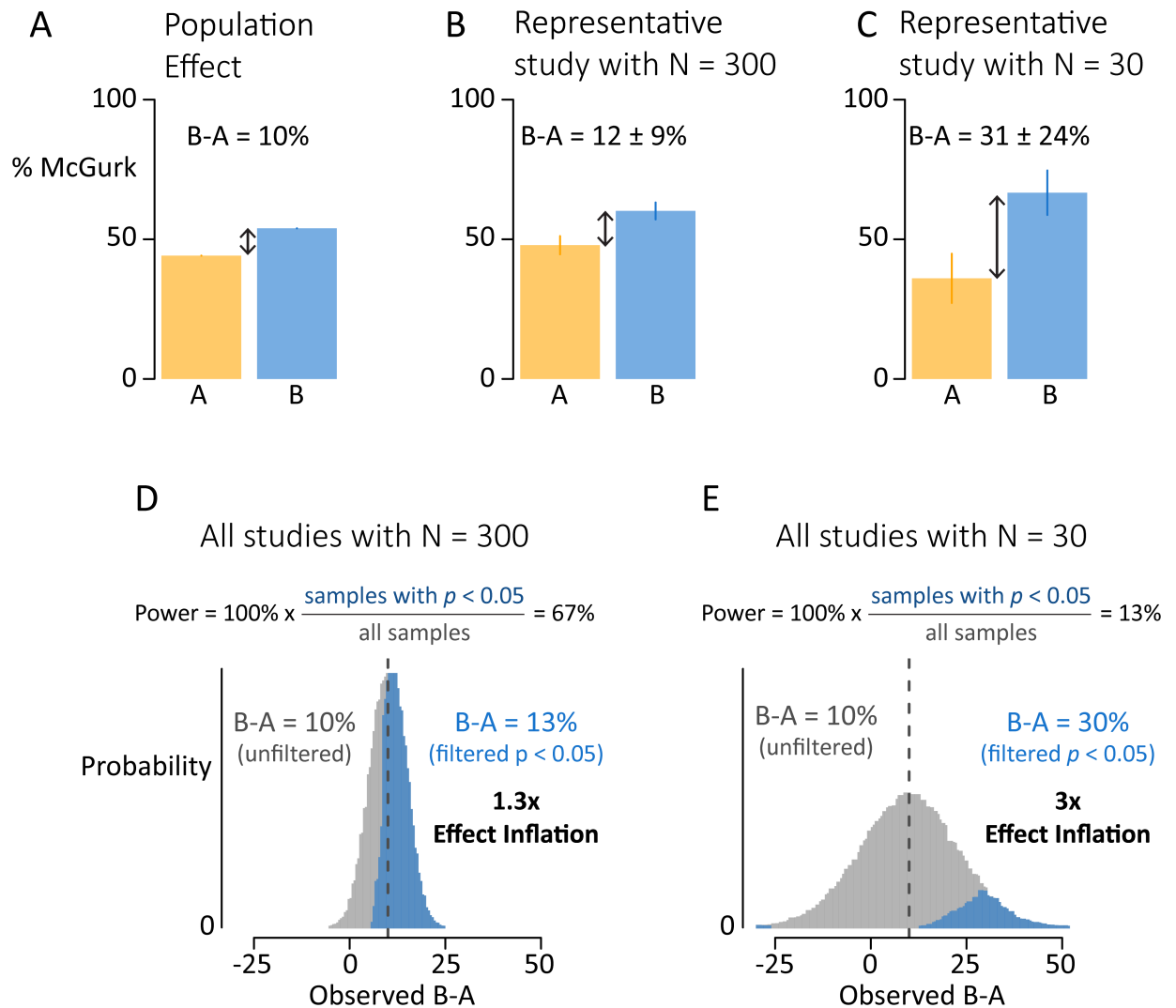


Fig 2. Small sample sizes lead to inflated estimates of population differences. (A) Two example populations with McGurk susceptibility of 45% (Population A, orange) and 55% (Population B, blue), producing a true difference in McGurk susceptibility of 10%. (B) A representative study using a sample of 300 subjects, 150 from each population. The estimated population difference after significance filtering is 12%, close to the true difference. (C) A representative study using a typical sample size of 30 subjects (15 from each population). The estimated population difference after significance filtering is 31%, far from the true difference. (D) Long run distribution of population difference estimates from studies with $N = 300$. Across all studies, the mean effect estimate is accurate (mean of 10%, gray and blue bars combined). Considering only studies with significant results (blue bars) inflates the mean estimate to 13% because only 67% of the studies are considered (power = 67%). (E) Long run distribution of population difference estimates from studies with $N = 30$. Considering only studies with significant results inflates the mean estimate to 30%, three times higher than the true difference on 10%. Across all studies, the mean effect estimate is accurate (mean of 10%, gray and blue bars combined).

<https://doi.org/10.1371/journal.pone.0202908.g002>

around the true difference; gray distribution in Fig 2E). This variance in effect estimation is large enough that in 3% of statistically significant studies, the population difference will be in the wrong direction: experiments will incorrectly conclude that McGurk susceptibility is higher in group A, rather than the true effect of higher susceptibility in group B. In contrast, across all $N = 300$ studies, every statistically significant study estimates the population difference in the correct direction.

Experimental manipulations to counter effect inflation

The first example demonstrated that a true difference between populations of 10% will be vastly overestimated using procedures that are common in the literature: sample sizes of 30 and a statistical significance filter. Since true effect sizes are rarely known in advance, and are often a motivator for performing the experiment, we next examined effect size inflation at a range of true population differences. As shown in Fig 3A, effect size inflation is most extreme

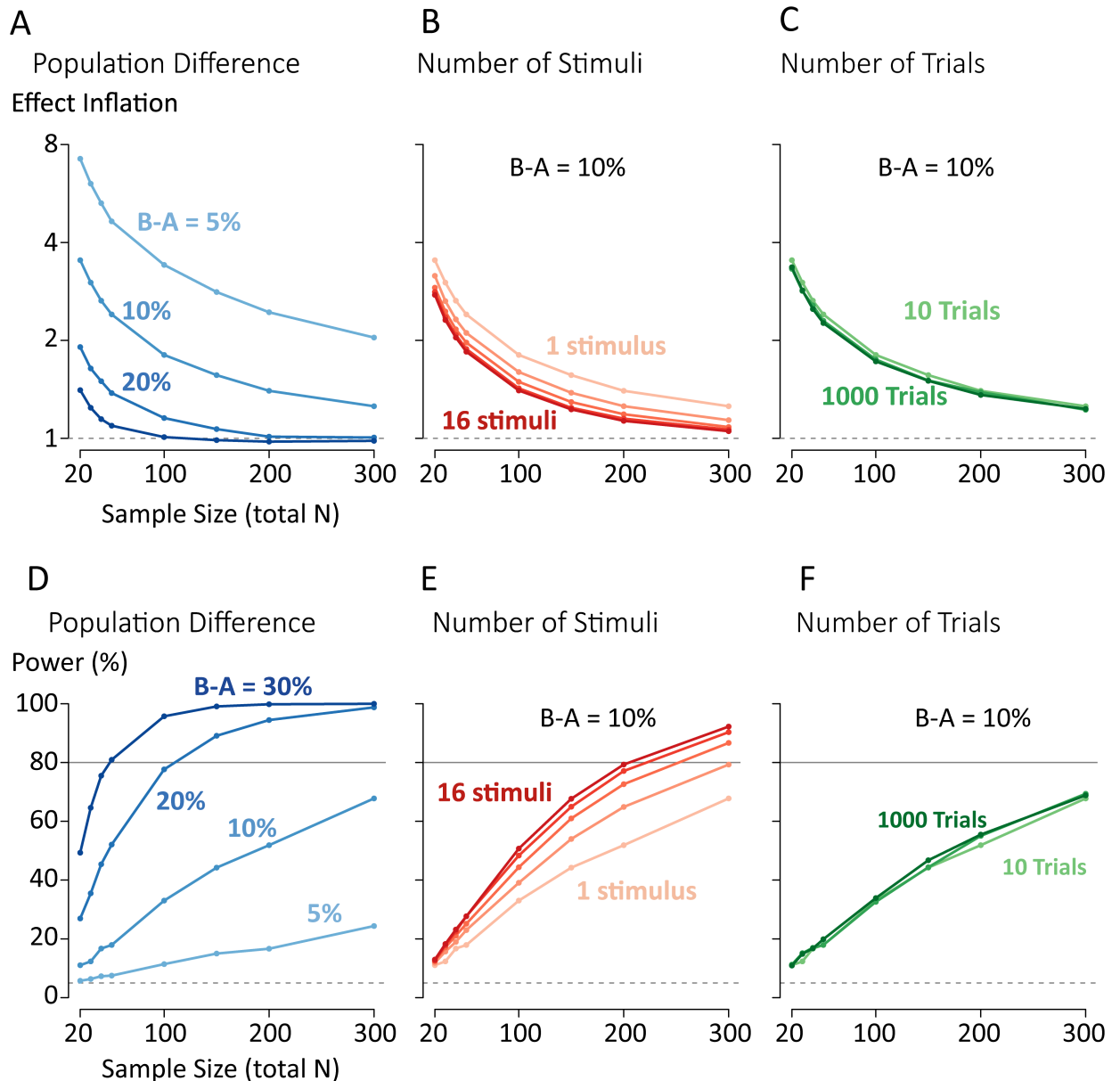


Fig 3. Factors contributing to effect inflation and statistical power in studies of population differences in McGurk susceptibility. In each panel, a statistical significance filter of $p < 0.05$ has been applied. (A) Effect inflation decreases as sample size increases for a fixed population difference. Doubling the known population difference (5%, 10%, 20%, 30%; blue lines) halves effect inflation at a given sample size. (B) Effect inflation decreases with increasing number of stimuli (1, 2, 4, 8, 16; red lines), shown for a fixed population difference of 10%. (C) Increasing the number of trials (10, 100, 1000; green lines) has little impact on effect inflation. (D) Statistical power (ability to detect a non-zero population difference) increases with increasing sample size and increasing population difference. (E) Increasing the number of stimuli produces diminishing returns for increasing statistical power. (F) Increasing the number of trials has minimal effect on statistical power.

<https://doi.org/10.1371/journal.pone.0202908.g003>

when population differences are small. For a true population difference of 5%, effect size inflation reaches over 7-fold for $N = 20$. In other words, on average, experiments with $N = 20$ (10 in each group) that select subjects from populations with true difference of 5% will publish effect estimates of over 35% (using $p < 0.05$ as the selection criterion).

Experiments can be designed to increase the true population effect and data can be collected from large numbers of subjects. Our simulations show the limitations of these approaches in isolation. A very large population difference of 20% cannot make up for a small sample size of 20: the resulting studies will still report a two-fold effect size inflation. Conversely, even a very large sample size of 300 cannot make up for a small population difference of 5% (effect size inflation of more than two-fold). Effect inflation occurs even at large sample sizes because of the common scientific practice of publishing only studies that show a statistically significant effect (commonly defined as $p < 0.05$). Removing statistical significance as a publication criterion prevents effect size inflation, but with small sample sizes, effect sizes will vary enormously from study to study, including possible changes in the direction of the effect.

What about other experimental manipulations? There is considerable procedural variation in the number of stimuli used to assess the McGurk effect, with some using a single stimulus, and others many more [21, 22]. To examine how stimulus count influences effect inflation, we modelled the effect of increasing stimuli for a true population difference of 10% (Fig 3B). In general, increasing the number of stimuli reduces effect inflation. The effect is most pronounced when increasing the number of stimuli from 1 to 2, which reduces the effect inflation from 1.8 to 1.6 at $N = 100$. Further increases in the number of stimuli produce diminishing returns, to an effect inflation of 1.2 for 16 stimuli.

Another experimental approach is to present only a single stimulus but increase the number of trials. A typical McGurk experiment may present 10 repetitions of a single stimulus to a subject. As shown in Fig 3C, even a vast increase in the number of trials, from 10 to 1000 results in little or no decrease in effect inflation, regardless of the number of subjects tested.

Relationship between effect inflation and statistical power

The statistical power of an experiment is defined as the probability that an experiment with a known true effect size and given sample size will result in a significant p -value. Power calculations implicitly incorporate a statistical significance filter, as they assume that the experimenter's goal is to detect a difference between populations with a significance of $p < 0.05$, rather than to accurately estimate the group difference in the long run. A power of 80% is often used as a benchmark in the literature [23], as it ensures that only one in five experiments will be discarded for not being significant (assuming that the true population difference is known).

As shown in Fig 3D, large sample sizes are required for adequate statistical power across the range of population differences. Even for a very large difference between-population difference of 20% an N of 110 is required to achieve 80% power. For a population difference of 5% and a sample size of N of 300, power is only 23%.

Increasing the number of stimuli can increase power (Fig 3E), but this increase is not enough to fully overcome small sample sizes or small population differences. For a population difference of 10%, the sample size required to achieve 80% power decreases from $N = 450$ with one stimulus, to $N = 300$ with 2 stimuli, to $N = 205$ with 16 stimuli.

Changing the number of trials has minimal effect on statistical power (Fig 3F). For a single stimulus and a population difference of 10%, 450 subjects are required to achieve 80% power with 10 trials; increasing to 1000 trials only reduces the required number of subjects to 405.

Discussion

Relevance to the published literature on the McGurk effect

Autism spectrum disorder (ASD) is frequently referred to as a disorder in which patients have impaired ability to integrate information across modalities [24]. This consensus is based on behavioral studies comparing people with ASD and healthy controls. Fig 4 shows a summary of 9 of these studies, with the difference in McGurk susceptibility between ASD and healthy controls plotted against the sample size of the studies. There is wide variation across studies, ranging from 45% less integration in ASD ($N = 34$; [25]) to 10% more integration in ASD ($N = 36$; [26]). These differences across studies have been attributed to interactions between clinical diagnosis and other factors, including stimulus, gender, temporal processing abilities, or participant age. For instance, one group reported a population difference for younger children but not older children [27] while another group reported the exact opposite effect [28].

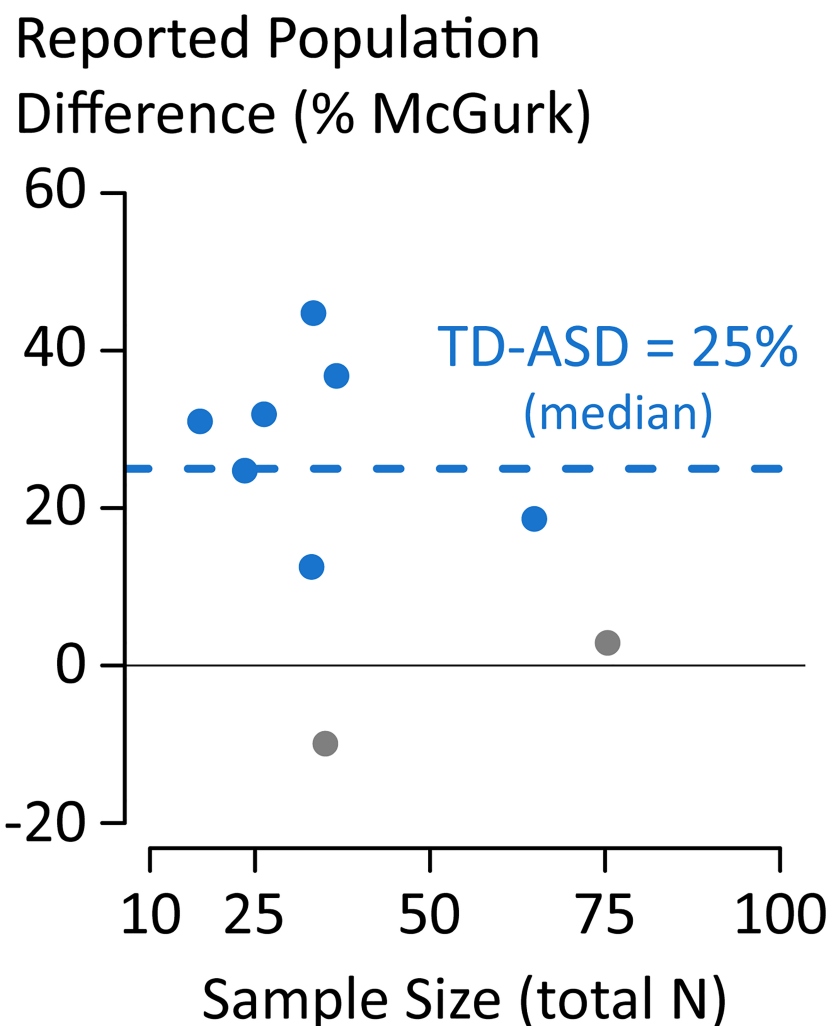


Fig 4. Reported population difference in McGurk susceptibility between children with autism spectrum disorder (ASD) and children with typical development (TD). Each symbol represents a single published study. Statistically significant ($p < 0.05$) differences are colored in blue, non-significant gray. Positive values indicate more McGurk susceptibility in the TD group vs. the ASD group. Median across studies is 25% (blue dashed line). Data for this table are available in the Supplemental Material.

<https://doi.org/10.1371/journal.pone.0202908.g004>

We suggest that a more likely explanation than a variety of *post hoc* moderators is that small sample sizes inevitably result in high variability in population effect estimates, causing wild swings in their magnitude and sign across studies [29, 30]. This contention is reinforced by noting that the ASD study with the largest sample size ($N = 76$) also reported the smallest effect, a population difference of only 3% [16]. Assuming that this large study is a more accurate measure of the true effect size, we can ask what the expected effect size will be for studies with a more typical sample size of $N = 36$ (mean across the 13 comparisons made in the ASD studies). Our simulations show that applying a significance-filter to studies of this size will inflate the population difference around nine-fold, resulting in an effect estimate of 28%, similar to the observed median of 25%.

A similar pattern is observed in studies of cultural differences in multisensory integration, which compare the prevalence of the McGurk effect between different linguistic or cultural populations. For instance, a number of studies have considered differences in McGurk susceptibility between speakers of tonal (*e.g.*, Mandarin) vs. non-tonal (*e.g.*, English) languages. Using small samples sizes, effect estimates have varied from +36% (*more* McGurk for non-tonal speakers; $N = 24$; [31]) to +17% ($N = 48$; [11]) to -8% (*less* McGurk for non-tonal speakers; $N = 40$; [32]). In contrast, a study with a large sample size found the smallest estimated difference: -4% (*less* McGurk for non-tonal speakers; $N = 376$; [15]).

Although these examples are taken from studies of population-level differences, the same problems arise in studies in which otherwise identical subjects receive different experimental manipulations, such as viewing different stimuli. In the studies of the McGurk effect we surveyed (Fig 1B), the median group size was 18 subjects. Because these studies included both between and within-group manipulations, the critical feature is the average *group size*, rather than the average total sample size.

Relevance to other studies of multisensory integration

Problematically small sample sizes in studies of multisensory integration are not restricted to examinations of the McGurk effect. Another common assay of multisensory integration is to measure a gain in performance when a unisensory cue is provided compared to a multisensory cue. For instance, one study found *reduced* multisensory gain for individuals with ASD ($N = 18$) compared to typically developing (TD) individuals ($N = 19$) when comparing auditory-only speech-in-noise perception with audiovisual speech-in-noise perception [33], while a separate lab found [34] found *increased* multisensory gain for individuals with ASD ($N = 16$) compared to TD individuals ($N = 14$) in a multisensory temporal order judgment. In line with the latter result, but at odds with the former, another study [35] suggested reported that individuals with ASD ($N = 29$) integrate over a larger temporal window than TD controls ($N = 17$). Although there are obvious task and stimulus differences between each of these studies, the ultimate goal of these studies (and those using the McGurk effect) is to establish generalizable estimates of group differences in multisensory integration. Critically, both multisensory speech-in-noise perception and multisensory temporal judgement tasks are known to have considerable variability even in healthy populations [36, 37]. For instance, in a population of 16 healthy controls, Magnotti and Beauchamp reported a range from 70 to 300 ms in sigma (a measure akin to the temporal binding window). Just as in the McGurk effect, large inter-individual variability in healthy controls makes the measure of intergroup differences with small sample sizes unreliable in general, and inaccurate (inflated) when only significant results are published.

Is this problem restricted to studies of multisensory integration?

Previous researchers have made formal arguments showing how using statistical significance as a publication filter distorts not just effect size estimates but also the perceived replicability of an effect, because only significant findings are published [38, 39]. Here we show how the presence of large individual differences exacerbates these problems. Because individual variation is a hallmark of many cognitive functions [40], large sample sizes are generally necessary to reliably measure population differences or experimental effects.

As with large individual differences, small sample sizes are also not unique to studies of multisensory integration: both the neuroscience [41] and psychology [42] literatures suffer from the same problem. Studies with small sample sizes over-estimate true effects, leading investigators in fruitless pursuit of the source of "large" effects despite failures to replicate [43]. While replication failures are often attributed to *post hoc* moderators or to the inevitable experimental differences between studies, our results show that a proliferation of small studies cannot resolve conflicting results. Instead, highly-powered studies aimed at accurate estimation are the only way to provide rigorous, reliable and reproducible studies of human behavior.

Suggestions for intergroup comparison studies of the McGurk effect

In this final section we summarize the results of our simulations as a series of suggestions for future inter-group comparisons of the McGurk effect.

Sample size. The major factor in determining the accuracy of the inter-group difference is the sample size. Assuming a true group difference of 20% studies with 100 subjects (50 in each group) have both reasonable statistical power (80%) and expected effect size inflation (1.2). For smaller expected effect sizes closer to 10%, 450 subjects are required for 80% power.

Stimulus number. In our simulations, increasing stimulus number had only moderate impacts on statistical power and effect size inflation. However, a complicating factor is that the stimuli themselves were highly variable in efficacy and the distribution of McGurk effect they elicited across subjects [13]. The choice of stimulus largely depends on the goal of the study. To study *group* differences in McGurk effect, pick a small number of relatively weak (or strong) stimulus with low variation in a control population. In contrast, to study *individual* differences, it is better to use a larger number of stimuli that show high variation.

Trial number. In our simulations, increasing the number of trials from 10 to 1000 had no effect on statistical power. The reason for this seemingly counterintuitive result is because of the higher variation *across* subjects than *within* subjects. Increasing the number of trials an individual is given will decrease the variability in our estimate of an individual's mean McGurk perception, but it has a lesser effect on our estimate of the group-level mean. An interesting extrapolation is how *few* trials could be used and still estimate group-level differences (although individual differences would not be accurately estimated). Using only 2 trials per participant (every participant will have 0%, 50%, or 100% McGurk), a sizeable population difference of 20% can be detected at 80% power with 150 participants, compared to 100 participants with 10 trials per participant (Fig 3). Whether such a tradeoff (more participants, fewer trials per participant) is worthwhile will depend on the particulars of the hypothesis being assessed.

Method

Bootstrap procedure to estimate effect inflation and statistical power

We used a bootstrapping procedure to create hypothetical replication datasets based on a large behavioral dataset ($N = 165$) collected in-person from Rice University undergraduate students

and described previously [13]. The goal of the simulation was to determine how experimental design choices impacted statistical power (ability to reject the null hypothesis) and effect estimation (accuracy of mean estimates from studies that reject the null hypothesis). We conducted the simulations using R [44]; source code is available on the authors' website: <http://www.openwetware.org/wiki/Beauchamp:DataSharing>.

The simulations proceeded in a series of 7 steps:

1. Simulation parameters were set. N : number of participants (20, 30, 40, 50, 100, 150, 200, or 300), E : size of the population difference (5%, 10%, 20%, and 30%), S : number of stimuli (1, 2, 4, 8 or 16), T : the number of trials for each stimulus (10, 100, or 1000).
2. McGurk perception rates, pM , for N participants at S stimuli were sampled with replacement from the empirical data and stored in a new dataset, D . This procedure ensures realistic within-subject correlations across stimuli.
3. Participants in D were randomly assigned to group A or B, ensuring equal group sizes (group size = $N / 2$).
4. To produce the population effect, the values of pM was shifted by E for all participants assigned to group B. Because of ceiling effects, the actual size of the shift needed to be greater than the desired mean population difference. We determined these values via simulation using sample sizes of 35,000 and stimulus count of 200. Actual shift values: 9.45%, 17.25%, 33%, and 55% produce population differences of 5%, 10%, 20%, and 30%, respectively. To ensure the simulation remained stochastic, pM was truncated to [5%, 95%] for all participants in D .
5. For each participant in D , we calculated an *observed* McGurk perception rate pF for each stimulus in S by sampling from a binomial distribution with T trials and true proportion pM .
6. A hypothesis test was conducted to compare the pF between groups A and B. For single-stimulus experiments, we used a two-sample, equal-variance t -test. For multi-stimulus experiment, we first averaged across stimuli within each subject, and then used a t -test. We used a t -test rather than a linear mixed-effects model here for computational efficiency. Because of how the population effect was created (a fixed shift for all participants and all stimuli), results for the LME and t -test were similar.
7. To obtain long-run behavior, we repeated these steps 35000 times for each parameter combination.

To summarize the results of the simulations we used two summary measures: Statistical Power: the proportion of simulations that rejected the null hypothesis, and Effect Inflation: the mean ratio of the estimated effect magnitude (absolute value of the mean difference between groups) and the true effect magnitude, for the subset of the simulations that rejected the null hypothesis, also called the expected Type M error or exaggeration ratio [29, 30]. Effect inflation provides a quantitative measure of the impact of the statistical significance filter on population effect estimates.

Summary of population differences in children with autism spectrum disorder

To assess the impact of small sample sizes in an important area of multisensory integration research, we reviewed studies that compared McGurk susceptibility between children with

autism spectrum disorder (ASD) and children with typical development (TD). We used Google scholar in the Fall of 2017 to find experimental studies comparing McGurk perception between individuals with ASD and controls. We looked for articles that cited the initial McGurk and MacDonald paper describing the illusion, using the keywords "autism", "ASD", and "Asperger's". We included only experimental studies (rather than reviews) to avoid including a specific dataset more than once. We included all studies for which group means and sample sizes were available. For studies with multiple group-level comparisons, we present them as separate data points. Two studies reported significant interactions but not group mean differences and are not included. A list of the studies used and the data collected are available from the authors' website: openwetware.org/Beauchamp:DataSharing.

Supporting information

S1 Dataset. Data file for McGurk studies comparing individuals with ASD and controls.

This file contains data from 10 studies (13 total comparisons) that compared the proportion of McGurk effect between individuals with ASD and control individuals. For each between-group comparison, we recorded the sample size and McGurk prevalence for each group. (XLSX)

Acknowledgments

We thank Kristen Smith for help in reviewing studies of the McGurk effect.

Author Contributions

Conceptualization: John F. Magnotti, Michael S. Beauchamp.

Formal analysis: John F. Magnotti.

Funding acquisition: Michael S. Beauchamp.

Methodology: John F. Magnotti, Michael S. Beauchamp.

Software: John F. Magnotti.

Supervision: Michael S. Beauchamp.

Visualization: John F. Magnotti.

Writing – original draft: John F. Magnotti.

Writing – review & editing: John F. Magnotti, Michael S. Beauchamp.

References

1. Sumbly WH, Pollack I. Visual contribution to speech intelligibility in noise. *J Acoust Soc Am*. 1954; 26(2):212–5.
2. Grant KW, Seitz PF. The use of visible speech cues for improving auditory detection of spoken sentences. *J Acoust Soc Am*. 2000; 108(3 Pt 1):1197–208. Epub 2000/09/29. PMID: [11008820](https://pubmed.ncbi.nlm.nih.gov/11008820/).
3. Ma WJ, Zhou X, Ross LA, Foxe JJ, Parra LC. Lip-reading aids word recognition most in moderate noise: a Bayesian explanation using high-dimensional feature space. *PLoS One*. 2009; 4(3):e4638. Epub 2009/03/05. <https://doi.org/10.1371/journal.pone.0004638> PMID: [19259259](https://pubmed.ncbi.nlm.nih.gov/19259259/).
4. McGurk H, MacDonald J. Hearing lips and seeing voices. *Nature*. 1976; 264(5588):746–8. PMID: [1012311](https://pubmed.ncbi.nlm.nih.gov/1012311/).
5. Traunmüller H, Öhrström N. Audiovisual perception of openness and lip rounding in front vowels. *Journal of Phonetics*. 2007; 35(2):244–58.
6. de Gelder B, Vroomen J, Van der Heide L. Face recognition and lip-reading in autism. *European Journal of Cognitive Psychology*. 1991; 3(1):69–86.

7. Tremblay C, Champoux F, Voss P, Bacon BA, Lepore F, Theoret H. Speech and non-speech audio-visual illusions: a developmental study. *PLoS One*. 2007; 2(1):e742. Epub 2007/08/22. <https://doi.org/10.1371/journal.pone.0000742> PMID: 17710142.
8. de Gelder B, Vroomen J, Annen L, Masthof E, Hodiament P. Audio-visual integration in schizophrenia. *Schizophr Res*. 2003; 59(2–3):211–8. PMID: 12414077.
9. Romero YR, Keil J, Balz J, Niedeggen M, Gallinat J, Senkowski D. Alpha-band oscillations reflect altered multisensory processing of the McGurk illusion in schizophrenia. *Frontiers in human neuroscience*. 2016;10.
10. Traunmüller H. Factors affecting visual influence on heard vowel roundedness: Web experiments with Swedes and Turks. *FONETIK* 2009. 2009:166.
11. Burnham D, Lau S, editors. The effect of tonal information on auditory reliance in the McGurk effect. *AVSP'98 International Conference on Auditory-Visual Speech Processing*; 1998.
12. Sekiyama K, Tohkura Y. McGurk effect in non-English listeners: few visual effects for Japanese subjects hearing Japanese syllables of high auditory intelligibility. *J Acoust Soc Am*. 1991; 90(4 Pt 1): 1797–805. PMID: 1960275.
13. Basu Mallick D, Magnotti JF, Beauchamp MS. Variability and stability in the McGurk effect: contributions of participants, stimuli, time, and response type. *Psychonomic Bulletin & Review*. 2015:1–9.
14. Strand J, Cooperman A, Rowe J, Simenstad A. Individual Differences in Susceptibility to the McGurk Effect: Links With Lipreading and Detecting Audiovisual Incongruity. *Journal of Speech, Language, and Hearing Research*. 2014; 57(6):2322–31. https://doi.org/10.1044/2014_JSLHR-H-14-0059 PMID: 25296272
15. Magnotti JF, Basu Mallick D, Feng G, Zhou B, Zhou W, Beauchamp MS. Similar frequency of the McGurk effect in large samples of native Mandarin Chinese and American English speakers. *Exp Brain Res*. 2015; 233(9):2581–6. Epub 2015/06/05. <https://doi.org/10.1007/s00221-015-4324-7> PMID: 26041554.
16. Stevenson RA, Segers M, Ncube BL, Black KR, Bebko JM, Ferber S, et al. The cascading influence of multisensory processing on speech perception in autism. *Autism*. 2017:1362361317704413. Epub 2017/05/17. <https://doi.org/10.1177/1362361317704413> PMID: 28506185.
17. Ferguson CJ, Heene M. A vast graveyard of undead theories publication bias and psychological science's aversion to the null. *Perspectives on Psychological Science*. 2012; 7(6):555–61. <https://doi.org/10.1177/1745691612459059> PMID: 26168112
18. Sterling TD. Publication decisions and their possible effects on inferences drawn from tests of significance—or vice versa. *Journal of the American statistical association*. 1959; 54(285):30–4.
19. Rosenthal R. The file drawer problem and tolerance for null results. *Psychol Bull*. 1979; 86(3):638.
20. Gelman A, Weakliem D. Of beauty, sex and power: Too little attention has been paid to the statistical challenges in estimating small effects. *Am Sci*. 2009; 97(4):310–6.
21. Stropahl M, Schellhardt S, Debener S. McGurk stimuli for the investigation of multisensory integration in cochlear implant users: The Oldenburg Audio Visual Speech Stimuli (OLAVS). *Psychon Bull Rev*. 2017; 24(3):863–72. Epub 2016/08/27. <https://doi.org/10.3758/s13423-016-1148-9> PMID: 27562763.
22. Magnotti JF, Beauchamp MS. The noisy encoding of disparity model of the McGurk effect. *Psychonomic Bulletin & Review*. 2015; 22(3):701–9.
23. Cohen J. A power primer. *Psychol Bull*. 1992; 112(1):155–9. Epub 1992/07/01. PMID: 19565683.
24. Rosenberg A, Patterson JS, Angelaki DE. A computational perspective on autism. *Proc Natl Acad Sci U S A*. 2015; 112(30):9158–65. Epub 2015/07/15. <https://doi.org/10.1073/pnas.1510583112> PMID: 26170299.
25. Bebko JM, Schroeder JH, Weiss JA. The McGurk effect in children with autism and Asperger syndrome. *Autism Res*. 2014; 7(1):50–9. <https://doi.org/10.1002/aur.1343> PMID: 24136870.
26. Woynaroski TG, Kwakye LD, Foss-Feig JH, Stevenson RA, Stone WL, Wallace MT. Multisensory Speech Perception in Children with Autism Spectrum Disorders. *J Autism Dev Disord*. 2013. <https://doi.org/10.1007/s10803-013-1836-5> PMID: 23624833.
27. Taylor N, Isaac C, Milne E. A comparison of the development of audiovisual integration in children with autism spectrum disorders and typically developing children. *Journal of autism and developmental disorders*. 2010; 40(11):1403–11. <https://doi.org/10.1007/s10803-010-1000-4> PMID: 20354776
28. Stevenson RA, Siemann JK, Woynaroski TG, Schneider BC, Eberly HE, Camarata SM, et al. Brief report: Arrested development of audiovisual speech perception in autism spectrum disorders. *J Autism Dev Disord*. 2014; 44(6):1470–7. <https://doi.org/10.1007/s10803-013-1992-7> PMID: 24218241.

29. Gelman A, Carlin J. Beyond Power Calculations: Assessing Type S (Sign) and Type M (Magnitude) Errors. *Perspect Psychol Sci*. 2014; 9(6):641–51. Epub 2015/07/18. <https://doi.org/10.1177/1745691614551642> PMID: 26186114.
30. Gelman A, Tuerlinckx F. Type S error rates for classical and Bayesian single and multiple comparison procedures. *Computational Statistics*. 2000; 15(3):373–90.
31. Sekiyama K. Cultural and linguistic factors in audiovisual speech processing: the McGurk effect in Chinese subjects. *Percept Psychophys*. 1997; 59(1):73–80. PMID: 9038409.
32. Chen Y, Hazan V, editors. Developmental factor in auditory-visual speech perception—the McGurk effect in Mandarin-Chinese and English speakers. AVSP'07 International Conference on Auditory-Visual Speech Processing; 2007.
33. Smith EG, Bennetto L. Audiovisual speech integration and lipreading in autism. *J Child Psychol Psychiatry*. 2007; 48(8):813–21. Epub 2007/08/09. <https://doi.org/10.1111/j.1469-7610.2007.01766.x> PMID: 17683453.
34. Kwakye LD, Foss-Feig JH, Cascio CJ, Stone WL, Wallace MT. Altered auditory and multisensory temporal processing in autism spectrum disorders. *Frontiers in integrative neuroscience*. 2011; 4:129. <https://doi.org/10.3389/fnint.2010.00129> PMID: 21258617
35. Foss-Feig JH, Kwakye LD, Cascio CJ, Burnette CP, Kadivar H, Stone WL, et al. An extended multisensory temporal binding window in autism spectrum disorders. *Exp Brain Res*. 2010; 203(2):381–9. Epub 2010/04/15. <https://doi.org/10.1007/s00221-010-2240-4> PMID: 20390256.
36. Van Engen KJ, Xie Z, Chandrasekaran B. Audiovisual sentence recognition not predicted by susceptibility to the McGurk effect. *Attention, Perception, & Psychophysics*. 2017; 79(2):396–403.
37. Magnotti JF, Ma WJ, Beauchamp MS. Causal inference of asynchronous audiovisual speech. *Frontiers in Psychology*. 2013; 4:798. <https://doi.org/10.3389/fpsyg.2013.00798> PMID: 24294207
38. Francis G. Publication bias and the failure of replication in experimental psychology. *Psychon Bull Rev*. 2012; 19(6):975–91. Epub 2012/10/12. <https://doi.org/10.3758/s13423-012-0322-y> PMID: 23055145.
39. Ioannidis JP, Trikalinos TA. An exploratory test for an excess of significant findings. *Clin Trials*. 2007; 4(3):245–53. Epub 2007/08/24. <https://doi.org/10.1177/1740774507079441> PMID: 17715249.
40. Vogel EK, Awh E. How to exploit diversity for scientific gain: Using individual differences to constrain cognitive theory. *Current Directions in Psychological Science*. 2008; 17(2):171–6.
41. Button KS, Ioannidis JP, Mokrysz C, Nosek BA, Flint J, Robinson ES, et al. Power failure: why small sample size undermines the reliability of neuroscience. *Nat Rev Neurosci*. 2013; 14(5):365–76. Epub 2013/04/11. <https://doi.org/10.1038/nrn3475> PMID: 23571845.
42. Bakker M, van Dijk A, Wicherts JM. The Rules of the Game Called Psychological Science. *Perspect Psychol Sci*. 2012; 7(6):543–54. Epub 2012/11/01. <https://doi.org/10.1177/1745691612459060> PMID: 26168111.
43. Open Science C. PSYCHOLOGY. Estimating the reproducibility of psychological science. *Science*. 2015; 349(6251):aac4716. Epub 2015/09/01. <https://doi.org/10.1126/science.aac4716> PMID: 26315443.
44. R Core Team. R: A language and environment for statistical computing. Vienna, Austria: R Foundation for Statistical Computing; 2017.