

Uneven Bars? Looking for Environmental Microaggression Effects in NCAA Women's Gymnastics*

Tommy Morgan[†]

Seth Cannon[‡]

January 14, 2025

[Link to most recent version](#)

Abstract

We study whether environmental microaggressions, a type of racial microaggression, affect the performances of NCAA Division I women gymnasts of color when competing at any given university hosting women's gymnastics meets between 2015–2024. NCAA gymnastics provides an excellent setting to test for such an effect because scores are assigned individually but meets are hosted and attended by entire teams. We use the FairFace computer vision model to predict gymnasts' race from their official photographs, and a comprehensive dataset of women's gymnastics scores then allows us to use a difference-in-differences model with a broad set of fixed effects to isolate racial gaps in scoring by host university. Across two specifications – first, comparing Black and White gymnasts, and second, comparing White gymnasts to all other gymnasts – we find no convincing evidence that gymnasts' scores are affected by an environmental microaggression effect.

Keywords: women's gymnastics, college sports, racial microaggression theory

JEL Codes: J15, Z2, Z20

*The authors are grateful to Zach Flynn for his help in compiling the original data used in this study and also thank Joe Price and Jeff Denning for helpful direction in the early stages of this project. † Corresponding author. Email: tmorga39@vols.utk.edu. ‡ Email: sethbcannon@wharton.upenn.edu.

Introduction

According to Sue, Capodulipo, and colleagues (2007), *racial microaggressions* are “brief, everyday exchanges that send denigrating messages to people of color because they belong to a racial minority group.” In their seminal paper on the topic, the authors identify nine broad categories of racial microaggressions, each with its own theme. One particular category of interest was referred to by Sue et. al. as *environmental microaggressions*, referencing “[m]acro-level microaggressions which are more apparent on systemic and environmental levels.” Two examples of this type of microaggression supplied by the authors include “[a] college or university with buildings that are all named after White heterosexual upper class males” and “[t]elevision shows and movies that feature predominantly White people, without representation of people of color”. Given that the history of women’s artistic gymnastics in America is well-known to be sparse of women of color (Wamsley, 2023; Reid, 2024) and that some NCAA universities have complicated racial histories (for example, Brigham Young University (Bergara, 2013) and several SEC schools (Berry, 2004)), it may be the case that gymnasts of color participating in a historically unrepresentative sport experience environmental microaggressions at certain universities. If this is the case, these adverse pressures could manifest themselves as negative effects on those gymnasts’ performance.

NCAA women’s artistic gymnastics competitions provide several unique mechanical advantages to a study of environmental microaggression effects in addition to the historical relevance described above. First, meet winners are decided at the team level, but scores are assigned to routines at the individual level. This individual scoring element allows us to look for an individual-by-host effect (i.e. a gymnast of a certain race doing worse at a certain place) that would be unidentifiable or obfuscated in pure team scoring contexts. Second, the frequency and prevalence of meets across the country provide an exceptionally large sample size across a broad range of competitors. Third, like in other college sports, NCAA gymnastics meets are usually hosted by a specific university at a consistent venue;

relatively few are hosted at neutral sites, especially early in the regular season of competition. This allows us to examine each host university as a consistent environment in which environmental microaggression effects may present themselves.

In order to identify a performance effect based in a racial microaggression, we need to assign a race prediction to each gymnast in our dataset. The NCAA collects self-reported demographic data for all of its student-athletes, but these “ground-truth” self-perceptions of race are not available to the public at the individual level. Instead, we apply the FairFace machine learning model created by Kärkkäinen and Joo (2021) to a headshot photograph of each gymnast and assign them to one of seven race categories. This allows us to make a relatively unbiased and consistent assignment of race using a model with excellent out-of-sample performance.

We combine these predictions with a newly assembled dataset of all NCAA women’s gymnastics scores from meets occurring between 2015-2024 to examine 1) whether Black gymnasts experience a negative performance effect relative to their White peers and 2) whether White gymnasts experience a positive performance effect relative to their non-White peers at each of the 64 D1 universities that hosted an NCAA meet over that time period. For each university, we use a differences-in-differences approach to analyze the performance of NCAA gymnasts who are on **visiting** teams over first ten meets of any season(s) in which they performed at a meet hosted by that university at least once. Though we initially find significant differential racial gaps in scoring at nine universities, most of these gaps are in the expected direction, and none of them survive corrections for multiple hypothesis testing.

Background & Related Literature

NCAA Women’s Artistic Gymnastics

We begin with a description of how NCAA women’s gymnastics meets are scored that relies heavily on Grimsley and Wright (2019), a thorough explanation of scoring in

NCAA women's gymnastics written by experienced journalists¹. Our summary also shares points about the sensitivity of gymnastics scoring with the Gymnastics section of Meissner et al. (2021). While our summary does not cite these articles for specific points (as it could also be considered common knowledge about the sport, especially to fans), they were nevertheless very helpful to its creation.

In women's artistic gymnastics, a regular season meet is composed of four events: vault, uneven bars, balance beam, and floor exercise. Each performance is scored out of 10 by two to four judges whose independent scores are averaged to a final performance score. The typical regular season meet has four judges – two from in-region and two from out of region – each judging two events, with two judges per event. When there are more than two teams at a given meet, at least eight individual judges judge one event each (still with two judges per event) with no rotation between events. Importantly, within a given meet, the set of judges that score each event is constant. In each event, five or six gymnasts from each team perform; if six gymnasts perform, the lowest of those six scores is dropped when calculating the overall team score for that event. After all four events are complete, the team scores are summed to compute each team's final meet score.

At the NCAA level, scores are determined by two factors: the “start value” of the routine, which is the score a gymnast would receive by performing their prepared routine perfectly, and deductions taken from the start value for technical or execution errors. Though an individual score can range anywhere from zero to 10 in each event, routines are required by rule to have at least a 9.4 point start value, and they score below 8.0 very rarely².

Because the practical range of scores is small, tiny differences in average scores

¹ Elizabeth Grimsley is the founder and editor-in-chief of College Gym News. Rebecca Wright is the current CNN Politics Photo Editor and a former Photo Editor for The Red & Black, a news organization that covers the University of Georgia.

² In our final sample of scores, less than 1% of scores are lower than 8.0; the 5th percentile score is 8.9, the 25th percentile 9.575, and the median 9.75.

separate elite teams from great and decent teams. A perfect team score would be at least five 10.0 routines in each event, giving team event scores of 50.0 each and a final team score of 200.0. In reality, only the best teams approach that threshold by the end of a given season. A team would be considered elite if it has the potential to hit a 198.00 meet score, which is obtainable only with an average score of 9.9 from every gymnast in every event across the entire meet; great teams can hit a 197.00 meet score (a 9.85 average performance score); and good teams could reach a 196.00 meet score (a 9.8 average performance score). These are differences of 0.05 points on average per routine, so losing even one-hundredth of a point (0.01) on a routine could be substantially harmful to team success; this motivates our research, as even a small environmental effect on scores could be meaningful to NCAA competition.

The unique attributes of artistic gymnastics meets offer us several key advantages. First, scores are assigned to gymnasts on an individual basis. This allows us to use individual routine scores to look for the presence of an environmental effect that would manifest itself at the individual-by-host level, i.e. a gymnast of a certain race experiencing a drop in performance at a meet hosted by a particular institution. This scoring model makes this individual-by-host analysis straightforward and differentiates this paper from research on other NCAA team sports like basketball and football (as in Dix (2019, 2021a), for example) in which individual performance is not always easy to fully isolate from team performance.

Second, our sample size is very large. Previous research investigating behavioral effects using women's gymnastics has primarily focused on elite-level gymnastics competitions, which do not happen as frequently as NCAA meets. These papers most frequently deal with race-agnostic biases present in judges and competitors, finding effects attributable to difficulty bias (Rotthoff, 2020), overall ordering bias (Morgan and Rotthoff, 2014; Rotthoff, 2015; Joustra et al., 2020) and the superstar effect (Meissner et al., 2021) at the highest level of the sport. However, because there are so many fewer elite gymnasts

than NCAA gymnasts, these papers can only analyze the performances of hundreds of elite-level gymnasts, whereas our sample includes thousands of gymnasts performing over multiple years.

Third, like in other college sports, NCAA gymnastics meets are usually hosted by a specific university at a consistent venue; relatively few are hosted at neutral sites, especially early in the regular season of competition. This also means that institution-hosted NCAA meets differ from elite meets in the consistency of their environment, as elite meets are hosted at various international venues that are not necessarily fixed, with the Summer Olympics being a classic example. This consistency allows us to examine each host university as an environment in which environmental microaggression effects may present themselves.

Racial Microaggression Theory

Recent research that investigates the effects of environmental racial microaggressions at the college level is often focused on qualitative interviews or surveys of Black students' experiences at predominantly White institutions (PWIs) (Mills, 2020; Holliday and Squires, 2020). This observation is also generally true of literature in this field historically, as evidenced by the many hundreds of papers based on interviewing Black students attending PWIs published from 1965-2013 that are summarized in Willie and Cunnigen (1981), Sedlacek (1987), and Holliday and Squires (2020). Also relevant to research on racial-environmental effects at the college level is Dix's body of work on sports programs at historically Black colleges and universities (or HBCUs) in which he shows teams from HBCUs experiencing negative performance effects while competing against PWIs in football (Dix, 2017, 2021a), men's basketball (Dix, 2022a,b), women's basketball (Dix, 2019, 2020a, 2022b), baseball (Dix, 2020b), softball (Dix, 2021b), and volleyball (Dix, 2023).

Much research also exists on racial biases within the world of professional sports.

This research usually focuses on racial biases in referee/judge decisions (as in Price and Wolfers, 2010; Parsons et al., 2011; Gallo et al., 2012; Rotthoff, 2020; Eiserloh et al., 2020; and Pelechrinis, 2023) or in fan/commentator preferences (as in Andersen and La Croix, 1991; Preston and Szymanski, 2008; Reilly and Witt, 2011; Principe and van Ours, 2022; and Quansah et al., 2023). These studies use data from professional sports leagues in many sports and around the world to show that racial biases can affect sports teams and players both in competitive outcomes and perceived value. We contribute to this vein of research on race effects in sports by studying one of its subtypes (environmental microaggression effects) in a novel setting (NCAA gymnastics).

Of particular relevance to this paper is Caselli et al. (2023), in which the authors show that African players in a professional Italian soccer league improved their performance when COVID-19 prevented fans from attending their games. They argue that this effect stems from the absence of overtly racist fan behavior, which is common in that league. As in our analysis of gymnasts' performance, the authors evaluated individual-level performance scores (in this case, those scores assigned algorithmically to individual soccer players based on in-game contributions) in a generalized fixed effects model that allows them to control for player- and match-based fixed effects. They model the effects of the removal of racial aggressions towards players, while we analyze the introduction of gymnasts to a potential environmental microaggression. We add to what Caselli et. al. found for professional athletes by estimating site-specific racial scoring gaps for college-level athletes.

Empirical Strategy

Model

To be able to attribute causality to any estimate we produce, we must be reasonably sure that we have controlled for as many other factors that could influence a gymnast's score as possible. We suggest a simple model of a gymnast's score in any given

event that breaks down influential factors into four principal categories:

$$\text{score} = f(\text{ability, preparation, environment, event})$$

In this model, ability-related factors could include, for example, the genetic makeup of a given gymnast, the age at which they began training, and the set of skills they have the physical capacity to perform. Preparation-related factors might include the quality of the team and coaches surrounding a gymnast, the number of years a gymnast has been competing, the types of skills a gymnast chooses to practice, and the number of meets that have already occurred in a season. Environment-related factors could include the relative competency or biases of judges at a given meet; the altitude, location, quality, and name of the venue; the time of the meet; and so on. Event-related factors adjust the score for the nuances of each event; for example, it is possible to fall off of the uneven bars, but not the floor exercise. Finally, because there is bound to be stochastic error in any model of human behavior, we also include an error term in our model.

Having put forth this general framework, we further posit that gymnasts and their coaches behave as score maximizers. We assume that these agents aim to maximize the score a gymnast will receive in a given event by selecting routines of appropriate difficulty and practicing them adequately; this implies that a gymnast and her coach have perfect information about the ability- and preparation-related elements of the proposed scoring model, and can therefore select the optimal routine to perform and lineup to set. We also assume that they have no control over the environment-related factors of her scoring, so if there is an environmental effect on scoring unknown to these agents, they may not be able to respond optimally. Determining whether such an effect exists at this level is the focus of this paper.

Data

RoadtoNationals.com has been the official statistical and rankings website of the Women's NCAA Gymnastics program since the summer of 2015 (Fredericks and Fredericks, 2013). It has been used as an accessible source for NCAA scoring and team ranking data in existing literature, as in Xiao (2022), Van Dyke et al. (2020), and Law (2020). To our knowledge, our dataset is the first comprehensive pre-processed source for these scores, as the data is not readily available for download at its source. In total, our dataset includes 230,088 scores received by 4,720 gymnasts over all 3,580 meets across all three NCAA divisions over the 2015-2024 seasons. We make our full dataset of NCAA women's gymnastics scores and all code used for the analysis in this paper available for future use³.

In addition to collecting data on scores, we also need to assign a race to each gymnast. Since we do not have access to the *individual-level* data that each gymnast reports to the NCAA, we use the FairFace race prediction computer vision model created by Kärkkäinen and Joo (2021). In order to apply the model, we collected a headshot photograph of each individual gymnast in our dataset; these consist of official photos from their university team's website for the vast majority of gymnasts and comparable photos obtained from news articles or social media when photos were unavailable from official sources. After collecting the photos, we put them each through the FairFace prediction model to classify each gymnast into one of seven race categories: White, Black, Indian, East Asian, Southeast Asian, Middle Eastern, and Latina. We then align our categories with reported NCAA categories as far as possible, arriving at a final set of four aligned categories: White, Black, Latina, and Other.

We are aware that assignment of race based on visual features, even by computer algorithm, subjects us to the "eye of the beholder" problem discussed in Fort and Gill (2000). For this reason, we compare our classifications to the *aggregated* racial

³ These will be made available via ICPSR (with corresponding code in a GitHub repository) after publication. At that point, this footnote will be updated to reflect that fact.

demographics data provided by the NCAA (National Collegiate Athletic Association, 2018) in Table 1 below.

Table 1

Comparing predicted & self-reported racial demographics

Race	Sample Year									
	2015	2016	2017	2018	2019	2020	2021	2022	2023	2024
Panel A: Scorers in our sample (predicted race)										
Black	7.6% (88)	7.6% (89)	7.7% (89)	8.4% (100)	9.0% (107)	9.8% (114)	9.6% (89)	9.4% (119)	10.1% (128)	10.6% (142)
Latina	4.4% (51)	5.6% (65)	6.2% (72)	6.7% (80)	7.1% (84)	7.2% (84)	8.2% (76)	9.0% (114)	9.6% (122)	8.5% (114)
Other	11.21% (130)	11.3% (132)	11.3% (131)	12.1% (144)	12.9% (153)	12.0% (140)	12.0% (112)	13.3% (169)	12.8% (163)	11.5% (154)
White	76.8% (891)	75.5% (882)	74.9% (870)	72.8% (865)	71.0% (843)	71.0% (828)	70.3% (654)	68.4% (868)	67.5% (857)	69.4% (931)
Unique Gymnasts	1,160	1,168	1,162	1,189	1,187	1,166	931	1,270	1,270	1,341
Panel B: All enrolled gymnasts (self-reported race)										
Black	7.8% (116)	7.9% (118)	7.7% (117)	8.0% (123)	8.2% (127)	8.5% (129)	8.4% (133)	7.7% (131)	8.1% (139)	7.7% (137)
Latina	3.9% (58)	4.4% (66)	4.3% (66)	3.8% (59)	3.2% (50)	3.6% (55)	4.0% (64)	4.6% (78)	5.3% (91)	6.2% (111)
Other	15.7% (234)	16.4% (246)	17.9% (273)	18.0% (279)	20.4% (315)	20.4% (309)	19.9% (316)	23.8% (407)	22.9% (392)	21.8% (390)
White	72.7% (1,084)	71.4% (1,072)	70.1% (1,067)	70.2% (1,086)	68.1% (1,050)	67.5% (1,022)	67.7% (1,073)	63.9% (1,091)	63.7% (1,093)	64.3% (1,147)
Unique Gymnasts	1,492	1,502	1,523	1,547	1,542	1,515	1,586	1,707	1,715	1,785

Note: Percentages of the total count of unique gymnasts from a given year with a given predicted/self-reported race are reported, with actual counts in parentheses.

It should be noted that the NCAA database includes all registered student-athlete gymnasts, whereas our data only includes those who competed and received at least one score in a given year. Even with this caveat, FairFace predicts many more gymnasts in our sample as Latina than are reported in the NCAA database. This is likely because FairFace only identifies gymnasts as a single race when they may identify in the NCAA demographics as Two or More Races or Unknown; this is especially likely to complicate the counts for the Latina category due to the complicated race vs. ethnicity issue that applies to them. The comparisons in Table 1, especially our overprediction of the Latina category,

motivate our decision to estimate our model in only two specifications: 1) comparing Black gymnasts to White gymnasts, looking for an environmental microaggression effect; and 2) comparing White gymnasts to all other (i.e. not White) gymnasts, looking for a sort of environmental micro-privilege effect.

Sample Construction

We begin by narrowing our sample to a subset of scores that meet four criteria: First, we drop scores from meets hosted at neutral sites (i.e. without a specific host university). Second, we drop meets with special titles like “SLC Regional”, “John Zuerlein Invite”, and “Big 12 Championships”. Third, we remove scores received by gymnasts who are not competing for Division I (DI) schools. Finally, in our Black-White comparisons, we drop scores from gymnasts that FairFace predicts are not Black or White (in our White-not White comparisons, we drop no scores for this step). But why narrow the sample in such a way?

First, we exclude any meet that is hosted at a neutral site in order to make the context surrounding the scores in our sample as similar as possible. Neutral site meets are hosted away from any university’s usual venue(s) by definition; therefore, they create a different type of environment than a typical regular season meet would have. This would confound the across-university comparison that is the focus of our research question, so we remove these meets.

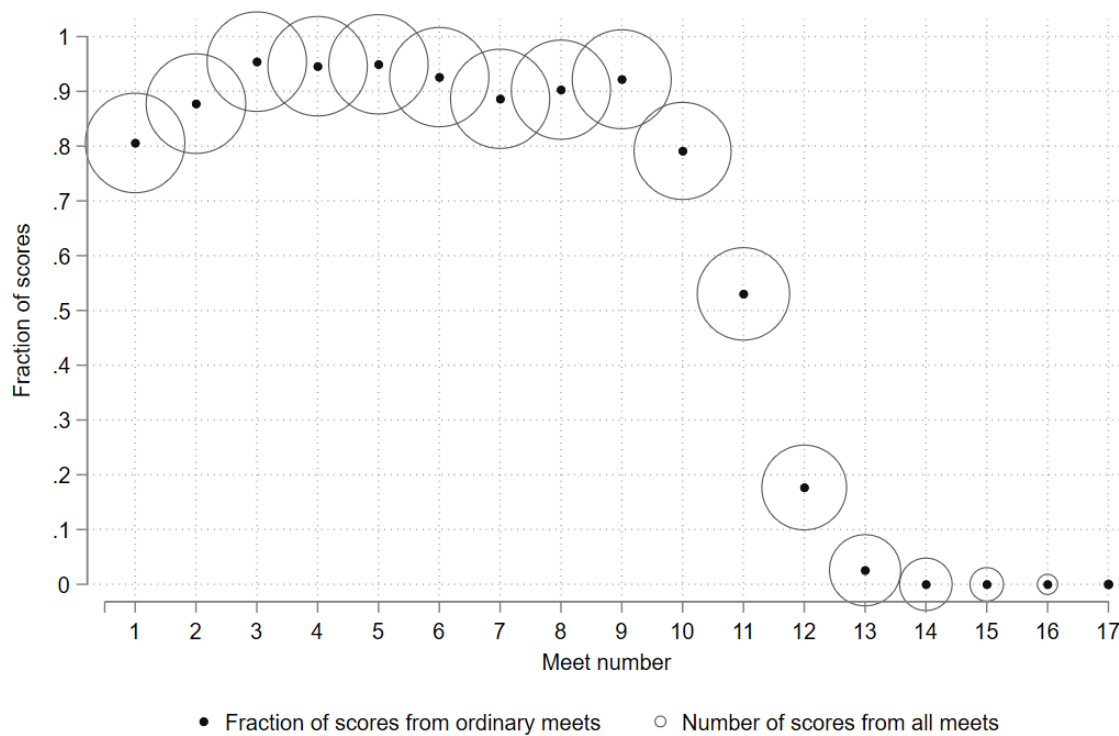
Second, we exclude titled meets in order to completely exclude playoff and invitational meets. Playoff meets can change incentives for gymnasts from the usual score maximization incentive⁴, and they may also present a higher pressure environment than

⁴ Early playoff meets advance the top teams from the meet to the next round, and final round playoff meets are won by teams, not gymnasts. A gymnast’s incentive to maximize her score may change if, for example, she is the sixth to compete on the uneven bars and she knows her team will qualify for the next round as long as her score is 9.75 or higher. In that case, she might adjust her planned routine to make it less likely she commits a major error instead of pushing for her highest possible score; this is the exact change in incentives we avoid by excluding playoff meets from our sample.

regular meets, representing another change that could muddy our central across-university comparisons. Invitational meets are also not conducive to our comparisons, as they are frequently hosted by organizations, not universities, which is a further change to the regular meet environment we wish to exclude.

Figure 1

Fraction of scores from “regular” meets by meet number.



Note: A “regular” meet is a meet hosted by a university without a special title, such as “invitational” or “regionals”. No meets hosted at neutral sites are included in calculations for this figure.

Figure 1 (pictured above) shows the fraction of the total set of scores that are from meets that survive our first and second sample narrowing criteria, with outer rings showing the relative number of scores at that meet number. As seen in the figure, about 80% of teams’ first meets are ordinary meets whose scores remain in our sample, whereas no team has a singular regular meet beyond their 13th meet. Since we drop mostly later meets with fewer scores under these criteria, Figure 1 shows that we keep the vast majority of meets in

our sample universe, especially those earlier in the season that are more likely to be ordinary meets.

After applying the first two sample narrowing criteria, we further limit our sample to scores received by DI gymnasts. Table 2 below illustrates average scoring by event, division, and race, and it shows that DI gymnasts are much better and more consistent scorers both within and across events than their non-DI counterparts.

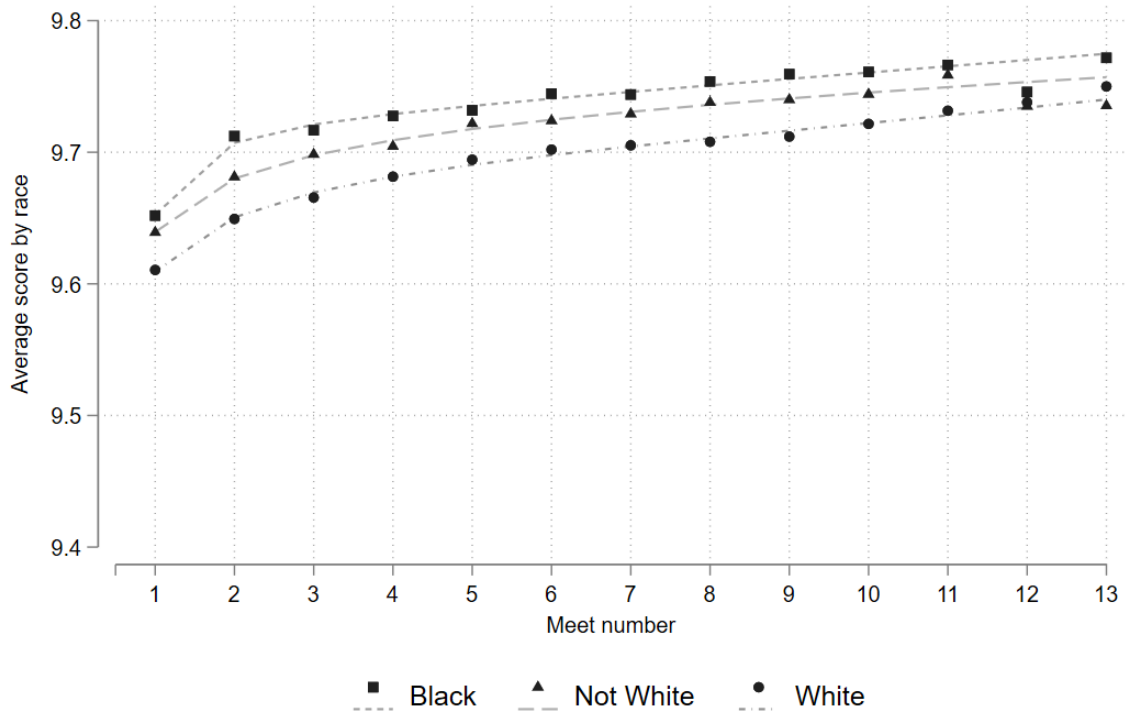
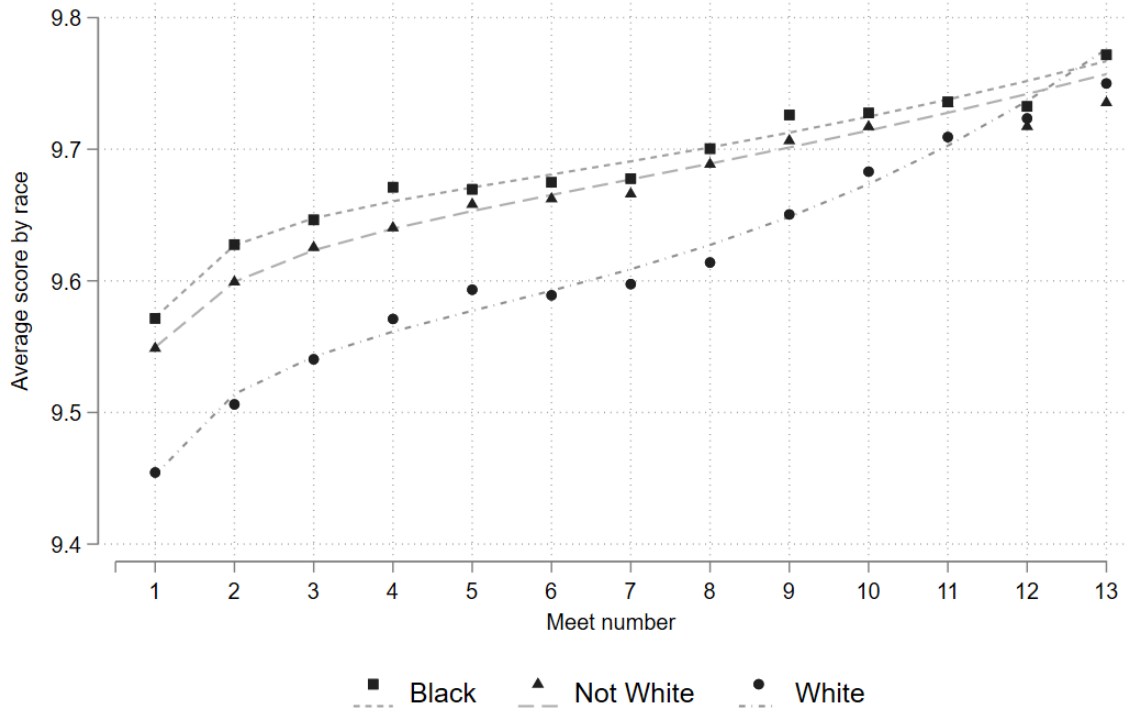
Table 2

Average score by event, NCAA division, and predicted race

	Vault	Bars	Beam	Floor	Overall
Panel 1: All gymnasts					
White	9.64 (0.26)	9.52 (0.50)	9.55 (0.41)	9.63 (0.37)	9.58 (0.40)
Black	9.71 (0.24)	9.62 (0.43)	9.61 (0.36)	9.72 (0.32)	9.67 (0.34)
All gymnasts	9.66 (0.25)	9.54 (0.48)	9.57 (0.39)	9.65 (0.36)	9.61 (0.38)
Panel 2: Only D1 gymnasts					
White	9.72 (0.18)	9.66 (0.36)	9.66 (0.30)	9.72 (0.30)	9.69 (0.30)
Black	9.76 (0.18)	9.70 (0.35)	9.68 (0.29)	9.77 (0.27)	9.73 (0.28)
All gymnasts	9.73 (0.18)	9.67 (0.36)	9.67 (0.30)	9.73 (0.30)	9.70 (0.29)
Panel 3: Only non-D1 gymnasts					
White	9.39 (0.30)	9.09 (0.60)	9.22 (0.49)	9.38 (0.42)	9.27 (0.48)
Black	9.44 (0.29)	9.16 (0.55)	9.28 (0.46)	9.46 (0.45)	9.35 (0.45)
All gymnasts	9.39 (0.30)	9.10 (0.59)	9.23 (0.49)	9.38 (0.42)	9.28 (0.48)

Note: Scores from titled or neutral meets are not included in calculations for this table.

Figure 2
Average Score by Race and Meet Week Number



Note: Lines represent observation count-weighted fractional polynomials of best fit. Scores from titled or neutral meets are not included in this figure.

We also motivate our decision to include only DI gymnasts' scores with Figure 2, pictured above. The top panel of Figure 2 shows the trend in scoring when including all gymnasts in our sample across all three NCAA divisions (Divisions I, II, and III), and it clearly shows that White gymnasts catch up to their not-White and Black peers as the season progresses. However, when we limit the sample to only DI gymnasts as in the bottom panel of Figure 2, we see that gymnasts of all predicted races improve at nearly the same rate throughout the season. This same-rate improvement is the exact parallel trend that makes our differences-in-differences design possible, so we keeping only the DI scores properly enables our identification strategy.

After applying all of the above criteria, we further narrow the sample each time we estimate our model for a given university. Explaining which meets we include in our sample for each university is easiest via an example, so suppose we take the University of Alabama as a host university. The University of Arkansas women's gymnastics team performed at Alabama every other year beginning in 2016 through 2024, so we include in the University of Alabama sample every score from every Arkansas gymnast at every one of their non-invitational regular season meets from those five years (2016, 2018, 2020, 2022, and 2024) in our sample. In contrast, the Alabama sample only includes University of Denver scores from 2019, as this is the only season in which Denver visited Alabama over our 2015-2024 time span. We repeat this process for every team that visited Alabama over that time period, collecting scores from meets in seasons in which they visited Alabama to eventually build the full Alabama sample as depicted below in Table 3. We then estimate the models we describe below for Alabama and then repeat this sample building and estimation process until we have done so for each of the 64 DI universities to have hosted a meet in that time span.

Table 3*Team-seasons included in the University of Alabama sample*

Team	2015	2016	2017	2018	2019	2020	2021	2022	2023	2024
Arizona	X									
Arkansas		X		X		X		X		X
Auburn	X		X		X		X		X	
Boise State	X		X						X	
Bowling Green					X					
Denver					X					
Florida	X		X		X		X		X	
Georgia		X		X		X		X		X
Illinois										X
Iowa State			X							
Kentucky		X		X		X	X	X		X
LSU	X		X		X		X		X	
Michigan					X					
Michigan State									X	
Minnesota										X
Missouri		X		X		X		X		X
North Carolina				X				X		
Northern Illinois					X					
Oklahoma	X			X		X				
S.E. Missouri					X					
Temple					X					
West Virginia		X								
Western Michigan								X		
Scores by Year	1,414	1,173	1,078	1,441	2,104	1,079	768	1,318	1,127	1,390
Total Sample Size	12,892									

Note: X denotes a year in which the scores from a team on a given row is included in the Alabama sample.

Estimation Strategy

We have two across-race comparison specifications: first, we compare Black gymnasts to White gymnasts, and then we compare White gymnasts to all other gymnasts. We begin both specifications by estimating a baseline differences-in-differences model. For

our Black gymnasts to White gymnasts model, it is:

$$\begin{aligned} \text{score}_{imet} = & \beta_0 + \beta_1 \text{Black}_i + \beta_2 \text{atHost}_m + \beta_3 \text{Black*atHost}_{im} \\ & + [\text{event}]_e + [\text{team}]_t + u_{iem} \end{aligned} \quad (1)$$

where subscripts i , m , e , and t refer to individual gymnasts, meets, events (vault, bars, beam, and floor), and teams, respectively. When we estimate our White-all others specification, we replace the Black indicator with an indicator for being White and adjust other variables accordingly. The dependent variable is the score earned by a gymnast in a given event, and the interaction term $(\text{Black*atHost})_{im}$ takes a value of 1 if the observation is a score received by a Black gymnast competing at a given host university and 0 otherwise. Black_i and atHost_m represent binary variables for a gymnast being Black and a meet being held at the focus university. We also include fixed effects for the event in which a given score was received denoted as $[\text{event}]_e$ (motivated by the differences in event means recorded in Table 2) and for the team a given gymnast is on denoted as $[\text{team}]_t$. We cluster our standard error u_{iem} at the event level. In this model, β_0 is the regression constant term, and the coefficient of interest is β_3 , which we interpret as the difference-in-differences estimate, or the differential impact of competing at the focus university on Black gymnasts relative to White gymnasts at that venue and those same Black gymnasts at other venues.

This baseline model certainly suffers from omitted variable bias. Chief among these omitted measures are 1) gymnast-specific characteristics such as build, talent, and training; and 2) meet-specific characteristics such as venue altitude, judge characteristics, and distance traveled for the meet. To control for several of these factors, we introduce gymnast and event-by-meet fixed effects:

$$\begin{aligned} \text{score}_{imet} = & \beta_0 + \beta_3 (\text{Black*atHost})_{im} + [\text{team}]_t \\ & + [\text{gymnast}]_i + [\text{event-by-meet}]_m + \epsilon_{iem} \end{aligned} \quad (2)$$

where each of the new square bracketed terms denotes a relevant set of fixed effects and the error term is represented by ϵ_{iem} , which we cluster at the event-by-meet level as explained below.

The gymnast fixed effects control for every event-, meet-, and time-invariant characteristic unique to a given gymnast, including her physical characteristics such as strength and height, her training before university, and her innate talent for gymnastics. The event-by-meet fixed effects – four indicators for each unique meet, i.e. $\mathbb{1}(\text{Vault at Utah State, 9 January 2015})$, $\mathbb{1}(\text{Balance Beam at Utah State, 9 January 2015})$, etc. – control for unchanging venue characteristics like geography, climate, and altitude, but they also control for judging. Within a given meet, the judges giving scores in a single event are constant: all scores in the floor exercise from the Utah State-hosted meet from 9 January 2015 are given by the same pair of judges. As a result, not only do the event-by-meet fixed effects capture event- and meet-constant characteristics, such as the date of the meet or the particular nuances of the floor exercise compared to the vault; they also control for the biases with which that pair of judges issued scores at that meet, because the judges giving the scores are an event-by-meet-constant characteristic.

For a causal interpretation of β_3 to be possible, we assume that Black (White) gymnasts performing at a given host university would experience the same relative change in performance at that venue as their White (not White) counterparts would in the absence of any environmental effect of that venue on Black (White) gymnasts' scoring. Since different teams visit different hosts at different points in the season throughout our dataset, we scrutinize our assumption with the bottom panel of Figure 2, in which we plot average scores in our sample by race and meet number. The figure shows that both Black, White, and non-White gymnasts see their scores increase over the course of a season on average, and that this trend is parallel across race groups. This figure demonstrates the baseline viability of the assumption needed for our estimates to be interpreted as causal.

We interpret a given host university's β_3 estimate as a gymnast-at-host-level effect.

If that effect is statistically significant at a given university, it could be due to an environmental microaggression factor like the name of a gymnasium or a predominantly non-Black student body, or it could be some other factor at a given university affecting Black gymnasts for some reason. If we were to see negative β_3 estimates within certain sets of universities – such as those included in a web article or Twitter thread compiling gymnasts speaking out against racism within their teams (such as in Duffy (2020) or Boswell (2020)) – then we may reasonably conclude that an environmental microaggression effect is present for Black gymnasts at some subset of universities. However, if universities with statistically significant β_3 estimates do not follow a noticeable trend, we might instead attribute those results to statistical noise, especially if those estimates do not survive corrections for multiple hypothesis testing.

Results

Pictured below, Table 4 lists the universities in our sample for which estimating Equations 1 and 2 returned estimates of β_3 significant at the 95% confidence level. Of the 64 DI host universities in each specification, only three ever have a significant Black-White coefficient, while six have significant White-all others coefficients. There does not appear to be a notable trend in what universities return these statistically significant estimates; they are not concentrated in any one geographic region, and most of them are not contained in Boswell’s compilation Twitter thread or Duffy’s article (Boswell, 2020; Duffy, 2020). In addition, only two estimates in the entire table (UC Davis and Towson) show an estimate with a sign opposite the direction of the gap that Table 2 shows already exists: White gymnasts are already the lowest average scorers, and Black gymnasts already the highest. Given these results, we find it unlikely that there is an environmental microaggression effect that affects Black gymnasts in the NCAA. But, then, how should we interpret the results we do see?

By using 95% confidence intervals as our judge for the statistical difference of β_3

Table 4*Universities with significant interaction term estimates*

Comparing Black gymnasts to White			Comparing White gymnasts to not White		
University	Estimate	<i>p</i> -value	University	Estimate	<i>p</i> -value
Sample size	(St. Err)		Sample size	(St. Err)	
Equation 1: Baseline diff-in-diff					
Auburn	0.025	0.017	Alaska	-0.066	0.017
N = 8,624	(0.005)		N = 2,822	(0.014)	
			Maryland	-0.014	0.017
			N = 16,741	(0.003)	
			Southern Utah	-0.031	0.038
			N = 13,563	(0.009)	
Equation 2: Full fixed effects model					
Pittsburgh	0.058	0.017	Alabama	-0.032	0.021
N = 11,643	(0.024)		N = 12,892	(0.014)	
UC Davis	-0.073*	0.032	Towson	0.042*	0.044
N = 8,304	(0.034)		N = 15,089	(0.021)	
			Washington	-0.036	0.030
			N = 11,338	(0.017)	

Note: Teams with significant estimates for each specification are included in alphabetical order.

**Sign of estimate indicates potential environmental effect (opposite of expected gap).*

from zero, we necessarily subject ourselves a 5% Type I error rate, meaning we expect to estimate a truly zero effect as statistically different from zero once in twenty tries. Since our testing is also two-sided, we would expect to estimate a truly zero effect as statistically greater than zero once in forty tries, and likewise in the opposite direction. If the null hypothesis of a true-zero effect held across all 64 universities in our sample, we would nonetheless expect to see one or two instances of statistically negative effects and one or two schools with statistically positive effects. In addition, because our sample already has an existing scoring gap, we might expect to see many more significant positive (negative) estimates in our Black-White (White-all others) specification, and fewer in the opposite

direction; indeed, this is the case, with only two of the nine significant estimates being in a direction opposite of the existing gap. As such, we argue that the few statistically significant estimates we observe in Table 4 are little more than what we would expect from random chance.

Given that we think our significant estimates are the result of random chance, we might ask a natural follow-up question: would a conservative adjustment to our standard errors designed to account for the multiple hypothesis tests we perform leave any school with a statistically significant estimate for β_3 ? After all, if significant results are the result of our method of statistical inference allowing too high of a Type I error rate (false rejection of the no effect null hypothesis), then a correction or adjustment designed to prevent Type 1 errors may help clarify our results. Following this logic, we apply the Bonferroni correction as established in Dunn (1961) and discussed in Armstrong (2014) and VanderWeele and Mathur (2019), in which a p -value is seen as significant only if it remains below 0.05 after being multiplied by the number of tests. Given that each model is estimated for 64 schools, any p -value that survives the correction would have to be smaller than $\frac{0.05}{64} = 0.00079$. Seeing that there are no p -values below 0.017 in Table 4 makes it clear that none of our estimates survive the correction, supporting our claim that the estimates are significant only by chance.

Discussion & Conclusion

We test whether Black (White) gymnasts experience a change in performance that their White (not White) competitors do not experience when competing at 64 Division I NCAA universities. Our comprehensive dataset of collegiate women's gymnastics scores allows us to isolate differential racial scoring gaps at a given host university by controlling for the influences of individual events, gymnasts, meets, and judges through a series of fixed effects. While we do initially find a few significant differences in score distributions using our model that could be attributed to environmental interaction effects, the fact that

none of them remain after adjusting for multiple hypothesis testing leads us to argue that our results provide no evidence that any given university hosting a given meet negatively affects Black (positively affects White) gymnasts' performance to a notable degree.

In short, we find no convincing evidence of an environmental effect at the universities in our sample for Black or White gymnasts. This result challenges expectations that systemic or environmental factors associated with universities with complicated racial histories might lead to measurable disparities in agents' performance at those universities. Instead, it may indicate that NCAA gymnastics environments, characterized by standardized scoring practices and strict competition protocols, limit the avenues through which such effects could manifest.

It may also be the case that DI gymnasts are simply good enough to overcome any potential environmental effect on their performance. We acknowledge that measuring the impact of environmental microaggressions only by examining performance scores is incomplete and overlooks the lived experiences the athletes might face by visiting these universities and stress that the absence of measurable performance effects does not eliminate the possibility that experiencing environmental microaggressions influences gymnasts in other ways. They may affect other dimensions of gymnasts' experiences, such as mental well-being or feelings of inclusion, which are beyond the scope of this study. Continuing to explore these aspects through qualitative methods could provide valuable context to complement our quantitative analysis.

Within this analysis, we also rely heavily on the FairFace computer vision model for race predictions, which imposes additional limitations. Though FairFace is known for its strong out-of-sample performance, it still relies on visual cues to assign race, which may not align with how individuals self-identify. With access to that self-reported data, future research could examine this question using self-identified race variables to better isolate the impact athletes who identify with given racial categories experience competing at different locations.

References

- Andersen, T. and La Croix, S. J. (1991). Customer racial discrimination in Major League Baseball. *Economic Inquiry*, 29(4):665–677. DOI: 10.1111/j.1465-7295.1991.tb00853.x.
- Armstrong, R. A. (2014). When to use the Bonferroni correction. *The Journal of the College of Optometrists*, 34:502–508. DOI: 10.1111/opo.12131.
- Bergara, G. J. (2013). "this time of crisis": The race-based anti-BYU athletic protests of 1968-1971. *Utah Historical Quarterly*, 81(2):204–229. DOI: 10.2307/45063320.
- Berry, M. S. (2004). Leveling the playing field: African-Americans and collegiate athletics. *EIU Historia*, 13:113–128. <https://www.eiu.edu/historia/Berry.pdf>.
- Boswell, S. F. (2020). Black gymnasts sharing their experiences of racism in the sport, a thread. *Twitter: @sf_boswell*. https://web.archive.org/web/20220128164941/https://twitter.com/sf_boswell/status/1270158886381764610.
- Caselli, M., Falco, P., and Mattera, G. (2023). When the stadium goes silent: How crowds affect the performance of discriminated groups. *Journal of Labor Economics*, 41(2):431–451. DOI: 10.1086/719967.
- Dix, A. (2017). A decade of referee bias against college football programs from Historically Black Colleges and Universities. *International Journal of Science Culture and Sport*, 5(3):197–212.
- Dix, A. (2019). "And 1" more piece of evidence of discrimination against Black basketball players. *Howard Journal of Communications*, 30(3):211–229. DOI: 10.1080/10646175.2018.1491434.
- Dix, A. (2020a). And 10 more years of bias against HBCU female basketball players. *Texas Speech Communication Journal*, 44:1–18.

- Dix, A. (2020b). Critical race theory, the NCAA, and college baseball: Contradiction on the diamond. In Milford, M. and Reichart Smith, L., editors, *Communication and Contradiction in the NCAA: An Unlevel Playing Field*, chapter 13, pages 213–234. Peter Lang, New York.
- Dix, A. (2021a). Referee judgments of communication in the field of play: A study on Historically Black Colleges and Universities in Division II college football. *International Journal of Sport Communication*, 14(4):554–573. DOI: 10.1123/ijsc.2021-0032.
- Dix, A. (2021b). Softball umpires call more walks than strikeouts when the pitcher plays for a Historically Black College and University. *International Journal of Sport Culture and Science*, 9(1):91–103. <https://web.archive.org/web/20231214222019/https://dergipark.org.tr/en/download/article-file/1412460>.
- Dix, A. (2022a). The non-Sweet Sixteen: Referee bias against Historically Black Colleges and Universities in men’s college basketball. *Sociology of Sport Journal*, 39(1):118–124. DOI: 10.1123/ssj.2020-0187.
- Dix, A. (2022b). Stay woke: An analysis of how referees evaluate the in-game communication of an HBCU that competes in a PWI conference. *Communication and Sport*, OnlineFirst. DOI: 10.1177/21674795221103407.
- Dix, A. (2023). Indications of referee bias in Division I women’s college volleyball: Testing expectancy violations and examining nonverbal communication. *International Journal of Sport Communication*, 16(4):414–422. DOI: 10.1123/ijsc.2023-0050.
- Duffy, P. (2020). Multiple NCAA gymnastics teams accused of racism by former athletes. *Gymnastics Now*. <https://web.archive.org/web/20240310002706/https://gymnastics-now.com/multiple-ncaa-gymnastics-programs-accused-of-racism>.
- Dunn, O. J. (1961). Multiple comparisons among means. *Journal of the American Statistical Association*, 56(293):52–64. DOI: 10.1080/01621459.1961.10482090.

- Eiserloh, D. G., Foreman, J. J., and Heintz, E. C. (2020). Racial bias in National Football League officiating. *Frontiers in Sociology*, 5(48). DOI: 10.3389/fsoc.2020.00048.
- Fort, R. and Gill, A. (2000). Race and ethnicity assessment in baseball card markets. *Journal of Sports Economics*, 1(1):21–38. DOI: 10.1177/152700250000100103.
- Fredericks, M. and Fredericks, J. (2013). Road to Nationals - NCAA gymnastics rankings. Accessed 12 June 2024 at <https://roadtonationals.com>.
- Gallo, E., Grund, T., and Reade, J. J. (2012). Punishing the foreigner: implicit discrimination in the Premier League based on oppositional identity. *Oxford Bulletin of Economics and Statistics*, 75(1):136–156. DOI: 10.1111/j.1468-0084.2012.00725.x.
- Grimsley, E. and Wright, R. (2019). Gymnastics 101: What to know about scoring, rankings and more before the next GymDog meet. *The Red and Black*.
https://web.archive.org/web/20230527011633/https://www.redandblack.com/sports/gymnastics-101-what-to-know-about-scoring-rankings-and-more-before-the-next-gymdog-meet/article_1008352e-56bd-11e2-b46e-0019bb30f31a.html.
- Holliday, N. R. and Squires, L. (2020). Sociolinguistic labor, linguistic climate, and race(ism) on campus: Black college students’ experiences with language at predominantly White institutions. *Journal of Sociolinguistics*, 25(3):418–437. DOI: 10.1111/josl.12438.
- Joustra, S. J., Koning, R. H., and Krumer, A. (2020). Order effects in elite gymnastics. *The Economist*, 169:21–35. DOI: 10.1007/s10645-020-09371-0.
- Kärkkäinen, K. and Joo, J. (2021). Fairface: Face attribute dataset for balanced race, gender, and age for bias measurement and mitigation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1548–1558.

- Law, J. (2020). A comparison of team effectiveness and perception of team success in academic and athletic teams. *Oregon State University Honors College Theses*.
- Meissner, L., Rai, A., and Rotthoff, K. W. (2021). The superstar effect in gymnastics. *Applied Economics*, 53(24):2791–2798. DOI: 10.1080/00036846.2020.1869170.
- Mills, K. J. (2020). “It’s systemic”: Environmental racial microaggressions experienced by Black undergraduates at a predominantly White institution. *Journal of Diversity in Higher Education*, 13(1):44–55. DOI: 10.1037/dhe0000121.
- Morgan, H. N. and Rotthoff, K. W. (2014). The harder the task, the higher the score: Findings of a difficulty bias. *Economic Inquiry*, 52(3):1014–1026. DOI: 10.1111/ecin.12074.
- National Collegiate Athletic Association (2018). NCAA demographics database. *NCAA.org*.
<https://www.ncaa.org/sports/2018/12/13/ncaa-demographics-database.aspx>. Accessed 15 December 2024, archived at <http://archive.today/K7QD4>.
- Parsons, C. A., Sulaeman, J., Yates, M. C., and Hamermesh, D. S. (2011). Strike three: Discrimination, incentives, and evaluation. *American Economic Review*, 101(4):1410–1435. DOI: 10.1257/aer.101.4.1410.
- Pelechrinis, K. (2023). Quantifying implicit biases in refereeing using NBA referees as a testbed. *Scientific Reports*, 13. DOI: 10.1038/s41598-023-31799-y.
- Preston, I. and Szymanski, S. (2008). Racial discrimination in English football. *Scottish Journal of Political Economy*, 47(4):342–363. DOI: 10.1111/1467-9485.00168.
- Price, J. and Wolfers, J. (2010). Racial discrimination among NBA referees. *The Quarterly Journal of Economics*, 125(4):1859–1887. DOI: 10.1162/qjec.2010.125.4.1859.

- Principe, F. and van Ours, J. C. (2022). Racial bias in newspaper ratings of professional football players. *European Economic Review*, 141. DOI: 10.1016/j.euroecorev.2021.103980.
- Quansah, T. K., Lang, M., and Frick, B. (2023). Color blind - Investigating customer-based discrimination in European soccer. *Current Issues in Sport Science*, 8(2):007. DOI: 10.36950/2023.2ciss007.
- Reid, T. (February 20, 2024). Making Black history in NCAA gymnastics. *College Gym News*.
<https://collegegymnews.com/2024/02/20/making-black-history-in-ncaa-gymnastics>.
- Reilly, B. and Witt, R. (2011). Disciplinary sanctions in English Premiership Football: Is there a racial dimension? *Labour Economics*, 18(3):360–370. DOI: 10.1016/j.labeco.2010.12.006.
- Rotthoff, K. W. (2015). (Not finding a) sequential order bias in elite level gymnastics. *Southern Economic Journal*, 81(3):724–741. DOI: 10.4284/0038-4038-2013.052.
- Rotthoff, K. W. (2020). Revisiting difficulty bias, and other forms of bias, in elite level gymnastics. *Journal of Sports Analytics*, 6(1):1–11. DOI: 10.3233/JSA-200272.
- Sedlacek, W. E. (1987). Black students on White campuses: 20 years of research. *Journal of College Student Personnel*, 28(6):484–495.
<https://psycnet.apa.org/record/1988-37333-001>.
- Van Dyke, E. D., Metzger, A., and Zizzi, S. J. (2020). Being mindful of perfectionism and performance among collegiate gymnasts: A person-centered approach. *Journal of Clinical Sport Psychology*, 15(2):143–161. DOI: 10.1123/jcsp.2019-0100.
- VanderWeele, T. J. and Mathur, M. B. (2019). Some desirable properties of the Bonferroni

correction: Is the Bonferroni correction really so bad? *American Journal of Epidemiology*, 188(3):617–618. DOI: 10.1093/aje/kwy250.

Wamsley, L. (September 2, 2023). Women’s gymnastics is changing in more ways than one. *NPR Sports*. <https://www.npr.org/2023/09/02/1197287064/womens-gymnastics-is-changing-in-more-ways-than-one>.

Willie, C. V. and Cunnigen, D. (1981). Black students in higher education: A review of studies, 1965-1980. *Annual Review of Sociology*, 7:177–198. <https://www.jstor.org/stable/2946027>.

Xiao, Y. J. (2022). Examining how region and individual competency affect team performance in Artistic Gymnastics using RStudio. *Scholarly Review Journal*, Fall 2022(4). DOI: 10.70121/001c.121677.