

Implicit Bias, “Science,” and Antidiscrimination Law

Samuel R. Bagenstos*

I. INTRODUCTION

In recent years, scholars of antidiscrimination law have increasingly focused on the problem of “implicit” or “unconscious” bias.¹ They have pointed to an expanding mass of evidence from experimental psychology that appears to demonstrate the pervasiveness of unconscious bias based on race, gender, and other legally protected characteristics, evidence that raises troubling questions about the effects of such bias on legally relevant behaviors. For example, research using such tools as the Implicit Association Test (IAT)—“which assesses bias by measuring the speed with which an individual associates a categorical status (such as black or white) with a given characteristic or description (such as good or bad)”²—has found that “[w]hite Americans, on average, show strong implicit preference for their own group and relative bias against African Americans.”³ Studies show that whites have similar biases against “other ethnic minority groups such as Latinos, Jews, Asians, and non-Americans,” as well as “the elderly, and women.”⁴ Interestingly, the studies also show that minorities and women often harbor the same implicit biases about their own groups that whites and men harbor against them.⁵

Legal scholars have used that psychological evidence to support a number of doctrinal proposals. For example, Linda Krieger and Susan Fiske have

* Professor of Law, Washington University School of Law. Thanks to Tristin Green, Rebecca Hollander-Blumoff, Jerry Kang, Laura Rosenbury, and, as always, Margo Schlanger for comments on an earlier draft of this Essay.

¹ For reviews of this literature, see Samuel R. Bagenstos, *The Structural Turn and the Limits of Antidiscrimination Law*, 94 CAL. L. REV. 1, 5–7 & nn.10–20 (2006); Jerry Kang, *Trojan Horses of Race*, 118 HARV. L. REV. 1489, 1497–1535 (2005). The implicit bias literature obviously builds on the classic works on unconscious bias and antidiscrimination law. See generally Linda Hamilton Krieger, *The Content of Our Categories: A Cognitive Bias Approach to Discrimination and Equal Employment Opportunity*, 47 STAN. L. REV. 1161 (1995); Charles R. Lawrence III, *The Id, the Ego, and Equal Protection: Reckoning with Unconscious Racism*, 39 STAN. L. REV. 317 (1987).

² Bagenstos, *supra* note 1, at 6. To take the IAT yourself, visit the IAT website, <https://implicit.harvard.edu/implicit/demo/takeatest.html>.

³ Nilanjana Dasgupta, *Implicit Ingroup Favoritism, Outgroup Favoritism, and Their Behavioral Manifestations*, 17 SOC. JUST. RES. 143, 147–48 (2004). For a good general overview of the research, see Anthony G. Greenwald & Linda Hamilton Krieger, *Implicit Bias: Scientific Foundations*, 94 CAL. L. REV. 945 (2006).

⁴ Dasgupta, *supra* note 3, at 147–48 (citations omitted).

⁵ See *id.* at 149.

used implicit bias research to criticize the “honest belief rule” and the “same actor inference”—two doctrines employed by many lower courts in addressing claims of intentional discrimination under Title VII of the Civil Rights Act of 1964.⁶ In addition, a number of scholars, most notably Jerry Kang and Mahzarin Banaji, and Christine Jolls and Cass Sunstein, have argued that psychologists’ implicit bias findings provide a firm justification for affirmative action programs in employment.⁷ Other scholars have argued that those findings justify reorienting employment discrimination law to take a more structural approach.⁸ Most controversially, Ian Ayres has suggested that IAT scores might “be used as a criterion for hiring both governmental and nongovernmental actors” to reduce the prevalence and effects of implicit bias.⁹

The legal academic literature has not been uniformly supportive of these recommendations. Writing before most of the evidence from the IAT had been reported, Amy Wax argued that “extending the framework created by existing antidiscrimination statutes to cover unconscious workplace disparate treatment is not a good idea because it is unlikely to serve the principal goals of a liability scheme—deterrence, compensation, insurance—in a cost effective manner.”¹⁰ In my previous piece on this issue, I agreed that the psychological findings present significant problems for the antidiscrimination project. But I argued that responding to implicit bias requires moving “beyond the generally accepted normative underpinnings of antidiscrimination law.”¹¹ As a result, I (despairingly) thought it unlikely that legislatures would adopt, and courts would carry out, effective responses to the problem.¹²

⁶ See Linda Hamilton Krieger & Susan T. Fiske, *Behavioral Realism in Employment Discrimination Law: Implicit Bias and Disparate Treatment*, 94 CAL. L. REV. 997, 1027–52 (2006). Under the “honest belief rule,” if an employer honestly believed that there was a nondiscriminatory reason for an adverse employment action, that belief can defeat a finding of discriminatory intent even if there was, in objective fact, no such reason. The “same actor inference” is the inference that if the person who fires an employee is the same one who hired that employee, the firing could not have been based on discriminatory intent.

⁷ See Christine Jolls & Cass R. Sunstein, *The Law of Implicit Bias*, 94 CAL. L. REV. 969, 978–88 (2006); Jerry Kang & Mahzarin R. Banaji, *Fair Measures: A Behavioral Realist Revision of “Affirmative Action,”* 94 CAL. L. REV. 1063 (2006). Michael Selmi made a similar point some years earlier. See Michael Selmi, *Testing for Equality: Merit, Efficiency, and the Affirmative Action Debate*, 42 UCLA L. REV. 1251, 1283 (1995).

⁸ See Tristin K. Green, *Discrimination in Workplace Dynamics: Toward a Structural Account of Disparate Treatment Theory*, 38 HARV. C.R.-C.L. L. REV. 91, 95–99 (2003); Susan Sturm, *Second Generation Employment Discrimination: A Structural Approach*, 101 COLUM. L. REV. 458, 460 & n.4 (2001). For a critique of structural proposals, see generally Bagenstos, *supra* note 1.

⁹ IAN AYRES, *PERVASIVE PREJUDICE? UNCONVENTIONAL EVIDENCE OF RACE AND GENDER DISCRIMINATION* 424–25 (2001).

¹⁰ Amy L. Wax, *Discrimination as Accident*, 74 IND. L.J. 1129, 1132–33 (1999).

¹¹ Bagenstos, *supra* note 1, at 3.

¹² *Id.* at 34–40. Linda Krieger, in her early work on unconscious bias, sounded similarly skeptical notes. See Krieger, *supra* note 1, at 1244–47.

Now, the arguments for using antidiscrimination law to respond to implicit bias face a new, more fundamental challenge. Gregory Mitchell and Philip Tetlock contend, in a recent piece, that the psychological research purporting to demonstrate the pervasiveness of implicit bias “fails to satisfy key scientific tests of validity.”¹³ Working through numerous studies of implicit bias, Mitchell and Tetlock contend that those studies fail tests of construct validity, internal validity, statistical validity, and external validity.¹⁴ Mitchell and Tetlock ultimately conclude that attaching the label of “science” to implicit bias research is “more honorific than descriptive,”¹⁵ and that the increasing acceptance of that research sets the stage “for an epistemic disaster of minor-epic proportions.”¹⁶

Mitchell and Tetlock make some effective points. In particular, they offer significant reasons for advocates of implicit bias research to be cautious before making sweeping statements that “science” compels one or another proposed legal reform or that only “hypocrisy and self-deception” could justify opposition to it.¹⁷ They also offer good reason to reject the suggestion that employers should use the Implicit Attitude Test to evaluate the implicit attitudes of candidates for managerial or hiring positions.¹⁸

But, as I hope to show in this Essay, Mitchell and Tetlock’s argument does not at all undermine the case for taking account of implicit bias in antidiscrimination policy. Even if one accepts every “scientific” critique they offer of the implicit bias literature—and there is substantial dispute within psychology on some of those critiques¹⁹—the case for using the law to respond to the problem of implicit bias remains strong. In the end, many of Mitchell and Tetlock’s critiques of implicit bias research rest, not on any scientific ground, but on normative assumptions about what kinds of

¹³ Gregory Mitchell & Philip E. Tetlock, *Antidiscrimination Law and the Perils of Mindreading*, 67 OHIO ST. L.J. 1023, 1023 (2006).

¹⁴ *See id.* at 1056–1115.

¹⁵ *Id.* at 1029.

¹⁶ *Id.* at 1118.

¹⁷ Mitchell and Tetlock single out a passage in which, they claim, Kang and Banaji do just that. *See id.* at 1029. The passage is: “If there is any value judgment embedded in behavioral realism besides those intrinsic to the scientific method, it is a second-order commitment against hypocrisy and self-deception. The law views itself as achieving just, fair, or at least reasonable results. If science reveals that the law is failing to do so because it is predicated on erroneous models of human behavior, then the law must transparently account for the gap instead of ignoring its existence.” Kang & Banaji, *supra* note 7, at 1065. I do not believe that Mitchell and Tetlock’s is the fairest interpretation of this passage. There, Kang and Banaji are making a general point about how the law ought to respond to science; they are not making any claim about what science shows. I certainly agree, however, that an appeal to “science” is a significant theme of Kang and Banaji’s article. In my view, appeals to “science,” of the type in which both Kang and Banaji and Mitchell and Tetlock engage, very often mask significant normative disputes. For an elaboration of the point in another antidiscrimination context, see Samuel R. Bagenstos, *The Americans with Disabilities Act as Risk Regulation*, 101 COLUM. L. REV. 1479 (2001).

¹⁸ *See* Mitchell & Tetlock, *supra* note 13, at 1115.

¹⁹ *See infra* note 23.

discrimination the law should seek to prevent and punish.²⁰ In particular, they rest on a very narrow view, based on notions of individual fault, that the law should prohibit only discrimination that results from self-conscious, irrational animus. That narrow view may resonate politically, but antidiscrimination doctrine and theory have consistently rejected it.²¹

Mitchell and Tetlock's argument is thus best understood, not as a scientific critique of implicit bias research, but as an argument about the normative bases for antidiscrimination law. That argument thus confirms my point that the major obstacle for advocates who seek to retool the law to address implicit bias is the widely accepted set of understandings about the goals of antidiscrimination law.²² Such advocates must therefore focus their efforts as much on developing the normative case for responding to implicit bias as on developing the scientific case that implicit bias exists.

My argument proceeds as follows. In Part II, I show that Mitchell and Tetlock's arguments do not call into question the case for using antidiscrimination law to respond to implicit bias. In Part III, I offer some thoughts about the normative challenges that remain for advocates of using antidiscrimination law in such a way. In Part IV, I present a brief conclusion.

II. THE RESILIENT CASE FOR USING ANTIDISCRIMINATION LAW TO RESPOND TO IMPLICIT BIAS

Mitchell and Tetlock amass what seems, cumulatively, to be an overwhelming phalanx of arguments against the science that lies behind proposals for using the law to combat implicit bias. When unpacked, however, nearly all of those arguments rest on a particular set of normative views about the kinds of bias to which antidiscrimination law ought properly to respond. Generally speaking, that view is a narrow one that treats discrimination as a wrong perpetrated by a discriminator who acts self-consciously and irrationally. But advocates of using the law to respond to implicit bias do not take that narrow view. To the contrary, they understand discrimination as a social problem that—whether or not it reflects the “fault” of any individual discriminator—has systematically harmful effects on the life chances of members of particular socially salient groups. Under that broader view of the problem of discrimination, the “scientific” evidence that Mitchell and Tetlock dismiss remains highly relevant and telling. Although they ostensibly attack the “science” of implicit bias research, Mitchell and Tet-

²⁰ In an earlier piece, Richard Banks and colleagues pointed out that criticisms of the race IAT often rest on normative disagreements about the nature of bias. See R. Richard Banks et al., *Discrimination and Implicit Bias in a Racially Unequal Society*, 94 CAL. L. REV. 1169, 1186–87 (2006). My argument goes further and pinpoints some of the ways in which Mitchell and Tetlock's argument rests on such normative disagreements.

²¹ For elaboration of the point, see generally Samuel R. Bagenstos, “*Rational Discrimination, Accommodation, and the Politics of (Disability) Civil Rights*,” 89 VA. L. REV. 825 (2003).

²² See Bagenstos, *supra* note 1, at 34–40.

lock's real target is the normative view of antidiscrimination law as reaching beyond acts reflecting the individual fault of the discriminator.

In this Part, I elaborate that point. In Section A, I show that many of Mitchell and Tetlock's criticisms of implicit bias research rest on the question-begging assumption that implicit bias is not a distinct phenomenon—that implicit bias cannot be real if it is not connected to overt and self-conscious prejudice. In Section B, I consider the possible non-“prejudice” explanations Mitchell and Tetlock offer for the studies' findings of implicit bias. Even if one accepts those explanations, I contend that they do not at all undermine the case for treating implicit bias as a serious social problem to which antidiscrimination law should respond. In Section C, by both reviewing the arguments addressed in the first two sections and introducing additional arguments made by Mitchell and Tetlock, I show that their objection to using the law to respond to implicit bias ultimately rests, not on any “scientific” objection, but on an individual-fault-based normative understanding of antidiscrimination law. This does not mean that the case for using the law to respond to implicit bias has been established. Mitchell and Tetlock point out some real questions that need to be answered in the implicit bias literature. And, given the significant consequences to individuals of using the IAT or other measures of implicit bias as a screening device for employment, their arguments should give employers pause before they use those measures in that way.²³ But it does mean that Mitchell and Tetlock have not even come close to delivering a fatal “scientific” blow to the implicit bias antidiscrimination program.

²³ For example, the link between implicit bias and discriminatory behavior that actually denies opportunities needs further development. See Mitchell & Tetlock, *supra* note 13, at 1067–69 (arguing that the studies show, at most, that implicit bias leads to “micro-level ‘discriminatory’ behaviors” that are *de minimis*). As workplaces evolve into more flexible, less hierarchical organizations, there is reason to believe that micro-level discriminatory behaviors will increasingly aggregate to cause denials of opportunities, see Bagenstos, *supra* note 1, at 11, but the “reason to believe” is not yet hard proof. For discussions of the proof that exists, see Greenwald & Krieger, *supra* note 3, at 953–55; Kang & Banaji, *supra* note 7, at 1072–75. The “‘shooter-bias’ studies,” see *infra* note 34, also provide evidence of a link between implicit bias and deprivation of opportunities. Mitchell & Tetlock, *supra* note 13, at 1068 n.146. At most, however, concern with the link between implicit bias and behavior calls for “future predictive validity studies” that address “the relationship between implicit associations and subjective evaluations of disadvantaged groups.” *Id.* at 1069. The relatively low correlation between the two most common methods of measuring implicit bias—the IAT and priming—also raises questions that need to be resolved, see *id.* at 1060–61, although even Mitchell and Tetlock acknowledge that the size of the correlation is the subject of debate within the psychological community, see *id.* at 1061 & n.128 (citing William A. Cunningham et al., *Implicit Attitude Measures: Consistency, Stability, and Convergent Validity*, 12 PSYCHOL. SCI. 163, 167 (2001)). Mitchell and Tetlock also raise important questions about the scaling of measures of reaction speed, see *id.* at 1091–93, and the importance of statistical outliers in the conclusions drawn by implicit bias studies, see *id.* at 1103–05. But these points call for more study of implicit bias measures, not for rejection of the implicit bias program.

A. *Begging the Question: Is Implicit Bias a Distinct Phenomenon?*

Many of Mitchell and Tetlock's criticisms of implicit bias research rest on the assumption that implicit bias reflects explicit prejudicial attitudes. Thus, they contend that implicit bias research lacks construct validity because the "empirical evidence" regarding the correlation between implicit bias and explicit prejudicial attitudes is "mixed."²⁴ Similarly, they contend that the implicit bias literature fails to "distinguish between automatically activated associations that, if called to people's attention, they would endorse and those associations that, if called to people's attention, they would categorically reject."²⁵ People may have negative, race-based implicit associations for a variety of reasons that do not stem from personal prejudice; those associations may instead reflect sympathy or awareness of cultural stereotypes and depressing realities. Such alternative explanations, Mitchell and Tetlock argue, "challenge the characterization of associations as tapping into racial attitudes—at least attitudes in the commonsense view that the attitudes imply an evaluative preference that, when brought to people's attention, they endorse and are even prepared to justify under appropriate conditions."²⁶

These arguments do not undermine the case for using law to respond to implicit bias. Scholars who advocate using the law in that way have strongly urged that implicit biases are meaningfully distinct from explicit attitudes.²⁷ Thus, "even people who express strongly egalitarian attitudes often show significant implicit biases."²⁸ And, even if those biases are called to their attention, the urge to rationalize them may be as strong as the desire to get rid of them.²⁹ To say that the concept of implicit bias lacks validity because implicit bias does not correlate empirically with explicit prejudice is therefore to assume the very conclusion that implicit bias scholars seek to challenge—that any "real" bias must be reflected in expressed attitudes (or the attitudes that would be expressed if the matter were called to one's attention).³⁰

²⁴ Mitchell & Tetlock, *supra* note 13, at 1062.

²⁵ *Id.* at 1084.

²⁶ *Id.* at 1080.

²⁷ See, e.g., Kang, *supra* note 1, at 1512.

²⁸ Bagenstos, *supra* note 1, at 7.

²⁹ See *id.* at 7–9.

³⁰ See Mahzarin R. Banaji et al., *No Place for Nostalgia in Science: A Response to Arkes and Tetlock*, 15 PSYCHOL. INQ. 279, 280 (2004) (responding to an earlier version of the argument made by Mitchell and Tetlock). Relying on an article by Blanton and Jaccard, Mitchell and Tetlock state that one of the advocates of implicit bias research, Mazharin Banaji, has taken two diametrically opposed positions on the question by treating a lack of correlation between implicit and explicit biases as evidence of validity in one article and then treating an observed correlation between the two as evidence of validity in another. See Mitchell & Tetlock, *supra* note 13, at 1062 n.130 (citing Hart Blanton & James Jaccard, *Arbitrary Metrics Redux*, 61 AM. PSYCHOL. 62, 66 (2006)); see also *id.* at 1095 (describing implicit prejudice advocates as "caught between a rock and a hard place" on this issue). It is not, I think, consistent with the implicit bias theory to treat a correlation with explicit prejudices

For the same reason, Mitchell and Tetlock's asserted "[d]ramatic disconfirmation[]" of the implicit bias theory—that "one study demonstrated that African-Americans actually preferred to interact with people classified as implicit racists"³¹—is neither dramatic nor a disconfirmation. A key point of implicit bias, to proponents of the theory, is that it is invisible.³² That minorities—many of whom may themselves share those biases³³—cannot detect implicit biases in others is not in any way a refutation of the theory. It would be a refutation only if one assumed the very point that implicit bias advocates seek to challenge—that such biases are visible.

B. A Narrow Theory of "Prejudice"

Relatedly, much of Mitchell and Tetlock's argument presumes that only a narrow sort of prejudice ought to be legally relevant. They suggest that, to the extent even blacks are shown to exhibit a bias against blacks, such bias can hardly be the sort of prejudice that ought to concern us. Thus, they dismiss the findings of "shooter-bias" studies, in which participants are more likely to shoot unarmed black targets and refrain from shooting armed white targets,³⁴ because black participants exhibited the same anti-black bias as did whites.³⁵ And they contend that "if the theoretical logic of the IAT requires that Jesse Jackson be classified as prejudiced against Blacks, then absurdity cannot be far off."³⁶

But that is hardly a scientific refutation of the implicit bias studies. As a logical matter, it is certainly possible that minority group members may be biased against members of their own group; it is hardly scientific to dismiss the possibility out of hand. Indeed, both experience and anti-discrimination doctrine make clear that the possibility is far from merely theoretical. Aside from such cultural clichés as the self-hating Jew,³⁷ antidiscrimination law and lore are full of stories of intentional discrimination by minority and female supervisors against members of their own race and sex.³⁸ Such acts of "self-hating discrimination" can have nearly as harm-

as validating the theory. But the major thrust of implicit bias research—including the overwhelming bulk of the research published by Banaji—treats the concept as one that is dissociated from explicit attitudes. At most, Mitchell and Tetlock's point is a minor "gotcha" that does not undermine the implicit bias theory.

³¹ *Id.* at 1065 (citing J. Nicole Shelton et al., *Ironic Effects of Racial Bias During Interracial Interactions*, 16 PSYCHOL. SCI. 397, 401 (2005)).

³² See Bagenstos, *supra* note 1, at 8 (discussing the argument by implicit bias scholars that such bias instantiates through "invisible cognitive processes").

³³ For further discussion of this point, see *infra* text accompanying notes 35–38.

³⁴ For discussion of these studies, see Kang, *supra* note 1, at 1525–26.

³⁵ See Mitchell & Tetlock, *supra* note 13, at 1068 n.146.

³⁶ *Id.* at 1085–86 n.204. Studies of implicit bias do show that minorities share the same biases as do whites, but to a substantially lesser degree than do whites. See Dasgupta, *supra* note 3, at 149.

³⁷ See generally SANDER L. GILMAN, *JEWISH SELF-HATRED: ANTI-SEMITISM AND THE HIDDEN LANGUAGE OF THE JEWS* (1986).

³⁸ For discussions in the law reviews, see Charles A. Sullivan, *Circling Back to the Ob-*

ful an effect on employment opportunities for minorities and women as can “other-hating discrimination” that occurs at the hands of white male supervisors. What should and does matter, for legal and policy purposes, is the fact that an individual is discriminated against; the race or gender identity of the discriminator is immaterial.

Mitchell and Tetlock also contend that “[t]here is strong evidence that psychological processes aside from out-group hostility can artificially inflate and otherwise distort scores on implicit measures such as the IAT.”³⁹ They then proceed to list a number of possible explanations for high implicit bias scores that do not reflect hostility toward minorities: figure-ground asymmetry (i.e., that “greater familiarity with one ethnic-racial group (e.g., Whites over Blacks) drives at least part of the race IAT effect”),⁴⁰ stereotype threat (i.e., that the fear of being labeled a bigot will drive some people to behave in a manner that appears bigoted),⁴¹ sympathy (i.e., the prospect that high IAT scores “are rooted more in compassion or guilt about the predicament of African Americans than in hostility or contempt”),⁴² and knowledge of the prejudice of others.⁴³ They offer similar explanations for the findings, in some implicit bias studies, that high IAT scores correlate with more awkward interactions with members of minority groups. Mitchell and Tetlock posit that shame about the way blacks have been treated in this country, or perhaps “social awkwardness stemming from lack of experience with members of other ethnic-racial

vios: The Convergence of Traditional and Reverse Discrimination in Title VII Proof, 46 WM. & MARY L. REV. 1031, 1085 (2004) (“In short, intra-racial or intra-gender discrimination occurs, if more rarely than cross-racial discrimination.”); Kenji Yoshino, *Assimilationist Bias in Equal Protection: The Visibility Presumption and the Case of “Don’t Ask, Don’t Tell,”* 108 YALE L.J. 485, 512 (1998) (“Self-hatred is a common response to [social] stigmatization, as exemplified by the instances of the self-hating black, Jew, or homosexual.”). See *Castaneda v. Partida*, 430 U.S. 482, 499 (1977) (“Because of the many facets of human motivation, it would be unwise to presume as a matter of law that human beings of one definable group will not discriminate against other members of their group.”); *Rine-smith v. Cent. County Fire & Rescue*, 156 Fed. Appx. 856, 859 (8th Cir. 2005) (“That two members of the Board of Directors were women does not preclude a finding of sex discrimination. Women may discriminate against women”); *Feingold v. New York*, 366 F.3d 138, 155 (2d Cir. 2004) (“We also reject the district court’s suggestion that an inference of discrimination cannot be drawn because Feingold was fired by another Jew.”); *Haywood v. Lucent Tech., Inc.*, 323 F.3d 524, 530 (7th Cir. 2003) (“Lucent begins by arguing that there should be a presumption of non-discrimination because Foote, who discharged Haywood, is also African-American. It is wrong; no such presumption exists, nor should one be created.”); *Ross v. Douglas County*, 234 F.3d 391, 396 (8th Cir. 2000) (“[W]e have no doubt that, as a matter of law, a black male could discriminate against another black male ‘because of such individual’s race.’”); cf. *Kadas v. MCI Systemhouse Corp.*, 255 F.3d 359, 361 (7th Cir. 2001) (finding it “altogether common and natural for older people . . . to be oblivious to the prejudices they hold, especially perhaps prejudices against the group to which they belong”).

³⁹ Mitchell & Tetlock, *supra* note 13, at 1072.

⁴⁰ *Id.* at 1075.

⁴¹ See *id.* at 1079–80.

⁴² *Id.* at 1081.

⁴³ See *id.* at 1084.

groups,” could be causing this effect.⁴⁴ “Of course,” they conclude, “a person experiencing such emotions and displaying awkward nonverbal behaviors is not necessarily prejudiced.”⁴⁵

But again, this is not a scientific disagreement so much as a normative one.⁴⁶ There is no reason why it should make a difference, from an anti-discrimination law and policy perspective, whether implicit bias reflects “hostility” toward minorities or any of the alternative explanations Mitchell and Tetlock offer. Whatever it reflects, implicit bias can result in behaviors and evaluations that limit the opportunities of minority group members. Racial unfamiliarity (“figure-ground asymmetry” in Mitchell and Tetlock’s technical term) exists in the employment setting as well as in the laboratory setting. To the extent that unfamiliarity with members of minority groups in certain work settings leads their white counterparts to engage (even unconsciously) in aversive behavior, that aversion will have the same negative effects on minorities’ opportunities as will aversion motivated by hostility.⁴⁷ To the extent that stereotype threat is a real risk in the laboratory setting, it can also operate in the employment setting. If whites act in a bigoted way on the IAT because they are afraid of being stereotyped as bigots, as Mitchell and Tetlock’s stereotype threat explanation suggests, whites who make employment decisions could just as well act in a bigoted way in those decisions because they are afraid of being stereotyped as bigots.

Sympathy, too, is hardly unproblematic. Paternalistic discrimination has a long history, and it has imposed substantial costs on members of disadvantaged groups. The classic statement of this point is Justice Brennan’s opinion in *Frontiero v. Richardson*,⁴⁸ which observed that sex discrimination was traditionally “rationalized by an attitude of ‘romantic paternalism’ which, in practical effect, put women, not on a pedestal, but in a cage.”⁴⁹ Moreover, work elaborating on the concept of “aversive racism”⁵⁰ shows that sympathy, guilt, social awkwardness, and shame have a major effect in reinforcing whites’ discrimination against blacks. As Linda Krieger summarizes these findings, “racial ambivalence, normative ambiguity, and

⁴⁴ *Id.* at 1096.

⁴⁵ *Id.* at 1097.

⁴⁶ See Banaji et al., *supra* note 30, at 282–83 (responding to an earlier version of Mitchell and Tetlock’s argument, without emphasizing the normative stakes).

⁴⁷ For this reason, I think Mitchell and Tetlock are unfair to Greewald and his colleagues. See Mitchell & Tetlock, *supra* note 13, at 1077 n.175.

⁴⁸ 411 U.S. 677 (1973).

⁴⁹ *Id.* at 684. Racial segregation in schools, too, was justified as being in the interests of blacks. See *Davis v. County Sch. Bd.*, 103 F. Supp. 337, 340 (E.D. Va. 1952), *rev’d sub nom.* *Brown v. Bd. of Educ.*, 347 U.S. 483 (1954); *Roberts v. City of Boston*, 59 Mass. 198, 209–10 (1849).

⁵⁰ See, e.g., John F. Dovidio & Samuel L. Gaertner, *On the Nature of Contemporary Prejudice: The Causes, Consequences, and Challenges of Aversive Racism*, in *CONFRONTING RACISM: THE PROBLEM AND THE RESPONSE* 3 (Jennifer L. Eberhardt & Susan T. Fiske eds., 1998).

fear of one's own potential prejudice all serve to amplify white's [sic] discrimination against blacks in the giving and requesting of assistance, the evaluation of behavior, physical distancing, and the selection of sanctions for social transgressions."⁵¹ If one cares about responding to conditions that operate as a systemic hindrance to the opportunities of members of minority groups—and that do so through no fault of the minority-group members themselves—one should care about responding to implicit bias, even if that bias is not rooted in hostility.

Finally, Mitchell and Tetlock contend that high implicit bias scores may not reflect prejudice but instead are “perfectly rational reactions to existing socioeconomic conditions.”⁵² In other words, because “African-Americans have, on average, fewer of the good things in life (high incomes and net worth, college educations, etc.) and more of the bad things in life (higher rates of imprisonment, violence, drug abuse, out of wedlock births, etc.),” it is perfectly rational, and indeed “inevitable,” that blackness will come to carry “a substantial network of negative associations.”⁵³

One might question just how rational these negative associations are. Are they precisely calibrated so that blackness carries negative freight only when and to the extent that the particular black person with whom one deals is statistically likely to be involved in “imprisonment, violence, drug abuse, out of wedlock births, etc.”?⁵⁴ Even if they are (which Mitchell and Tetlock make no effort to prove), that is entirely irrelevant from the perspective of antidiscrimination law and policy. Antidiscrimination laws do not exempt “rational” discrimination from their prohibitions. To the contrary, “[t]he prohibition of rational discrimination is a central component of anti-discrimination doctrine—and it may be the most important aspect of antidiscrimination law on the ground.”⁵⁵ It is well established, for example, that Title VII of the Civil Rights Act of 1964 prohibits employers from discriminating on the basis of race or sex even when it is demonstrably more costly to hire workers of a particular race or sex or to assign them to particular positions.⁵⁶ Title VII also prohibits employers from discriminating

⁵¹ Krieger, *supra* note 1, at 1240 (citations omitted).

⁵² Mitchell & Tetlock, *supra* note 13, at 1085.

⁵³ *Id.* at 1086.

⁵⁴ *Id.* Work on cognitive heuristics suggests that they will not be so precisely calibrated. See Cass R. Sunstein, *Probability Neglect: Emotions, Worst Cases, and Law*, 112 *YALE L.J.* 61, 70–83 (2002).

⁵⁵ Bagenstos, *supra* note 21, at 848; see also David A. Strauss, *The Myth of Color-blindness*, 1986 *SUP. CT. REV.* 99, 113–16 (1986).

⁵⁶ See *Int'l Union v. Johnson Controls, Inc.*, 499 U.S. 187, 210 (1991) (“The extra cost of employing members of one sex, however, does not provide an affirmative Title VII defense for a discriminatory refusal to hire members of that gender.”); *Los Angeles Dep't of Water & Power v. Manhart*, 435 U.S. 702, 716–17 (1978) (rejecting any “cost justification defense” under Title VII). Mitchell and Tetlock misstate the problem when they say that “the courts permit employers to base decisions on job-relevant attributes correlated with protected-category membership.” Mitchell & Tetlock, *supra* note 13, at 1087. Implicit bias, in their rational-association explanation, is more akin to basing decisions on protected-category membership because it is correlated with job-relevant attributes. But that is ex-

in order to satisfy customer preferences—even when customers will, in fact, deny their business to a firm that refuses to discriminate.⁵⁷ The law does so because even rational acts of discrimination can aggregate to cause systematic and cumulative disadvantage to members of minority groups.⁵⁸ Accordingly, even if implicit biases reflect purely rational discrimination, they remain extremely important for antidiscrimination law and policy. To disagree with that proposition reflects, not a scientific disagreement, but a normative judgment that rational discrimination is not the kind of discrimination the law should try to prohibit.

For the same reasons, Mitchell and Tetlock miss the mark when they assert that implicit bias cannot truly be prejudice because an individual's implicit bias score “depend[s] greatly on the specific stimuli chosen for the sorting task.”⁵⁹ They ask, rhetorically, “[i]f a relative aversion toward a category only materializes when we use certain types of names or faces—or when we embed those stimuli in certain types of context cues—are we justified in saying people are prejudiced toward the category as a whole or only toward particular exemplars of that category in certain contexts?”⁶⁰ The point of antidiscrimination law, however, is not to identify and punish prejudice as an inherent personal quality. The point is to prevent and provide remedies for conduct that deprives minorities of opportunities. Given that goal of antidiscrimination law, the fact that implicit bias may depend on the surrounding stimuli does not suggest that we should *ignore* that bias. Rather, it suggests that law and policy should aim to promote the inclusion in workplaces of the kind of stimuli that reduce implicit bias and its effects. And indeed, that is precisely the conclusion drawn by scholars who advocate using the law to respond to implicit bias.⁶¹ Again, the point is a normative one, not a scientific one.

actly what *Johnson Controls* and *Manhart* prohibit.

⁵⁷ See, e.g., *Fernandez v. Wynn Oil Co.*, 653 F.2d 1273, 1276–77 (9th Cir. 1981); *Diaz v. Pan Am. World Airways*, 442 F.2d 385, 389 (5th Cir. 1971).

⁵⁸ See Cass R. Sunstein, *The Anticaste Principle*, 92 MICH. L. REV. 2410, 2418 (1994) (“[T]he most elementary antidiscrimination principle singles out one kind of economically rational stereotyping and condemns it, on the theory that such stereotyping has the harmful long-term consequence of perpetuating group-based inequalities.”). For an extensive effort to justify the prohibition of rational discrimination, see Bagenstos, *supra* note 21, at 848–59.

⁵⁹ Mitchell & Tetlock, *supra* note 13, at 1106. For a similar point, see *id.* at 1113 (noting that photos of black persons “embedded in an egalitarian and positive-affect setting” do not trigger high implicit-bias scores).

⁶⁰ *Id.* at 1107.

⁶¹ See, e.g., Greenwald & Krieger, *supra* note 3, at 963–64; Christine Jolls, *Antidiscrimination Law's Effects on Implicit Bias*, in BEHAVIORAL ANALYSES OF WORKPLACE DISCRIMINATION (G. Mitu Gulati & Michael Yelnosky eds., forthcoming 2007); Kang & Banaji, *supra* note 7, at 1101–15. A similar point could be made about Mitchell and Tetlock's contention that implicit bias research lacks external validity because workplaces (at least those that embody “best-practices precepts in organizational behavior”) differ from experimental settings. See Mitchell & Tetlock, *supra* note 13, at 1107–15. External validity is always a difficult problem in experimental psychological research, and implicit bias research is not different. But if “best-practices precepts” can forestall discriminatory ef-

C. An Irrational-Animus Theory of Antidiscrimination Law

Much of the discussion above suggests that what really drives Mitchell and Tetlock's argument is not a scientific disagreement with the implicit bias researchers, so much as a normative disagreement about the appropriate scope and targets of antidiscrimination law. The rational-discrimination point may be especially telling: noting that "the IAT makes it possible to be a Bayesian bigot," they suggest (and cite an article in which Tetlock contends) that the possibility constitutes "a *reductio ad absurdum* of the IAT research program."⁶² But it is a *reductio* only if one takes the view that rational discrimination is not a proper target of antidiscrimination law, a point that most civil rights advocates would vigorously deny and that antidiscrimination doctrine categorically rejects.⁶³ This is not a debate about science; it is a debate about the proper scope of antidiscrimination law.

The assertion that rational discrimination does not count as discrimination suggests that Mitchell and Tetlock view matters through what Alan Freedman called the "perpetrator perspective"—the notion that what matters is whether the person accused of discrimination was at "fault," and not whether the person accused of discrimination actually caused or contributed to group-patterned harm.⁶⁴ That perspective is evident throughout the paper. Thus, when considering the possibility that implicit biases might be activated by sympathy, rather than antipathy, toward minority

fects of implicit bias, that point suggests that more workplaces should adopt those precepts. It is worth noting that many workplaces do not incorporate those precepts. When workplaces do, it is often a response to legal antidiscrimination mandates. *See, e.g., id.* at 1109 (noting that "hiring and staffing managers at many organizations," unlike laboratory subjects, "have been alerted to . . . the dangers of discrimination" and "are well aware of the need to have adequate and legal justifications for their judgments and decisions"); *id.* at 1110 (noting that managers, unlike laboratory subjects, "often receive clear and repeated admonishment" that they may not "allow job-performance-irrelevant characteristics, such as membership in ethnic-racial groups, to affect their personnel judgments"); *id.* at 1114 (noting "intergroup contact is associated with reduced prejudice, including 'prejudice' at the implicit level") (footnote omitted). To the extent that legal rules encourage adoption of these best-practices precepts—and therefore promote intergroup contact at the workplace—that is an argument *for*, not against, legal intervention in those workplaces that have not adopted them. Relatedly, Mitchell and Tetlock's argument that the "stranger-stranger interactions employed in" implicit bias experiments "do not reflect the types of interactions that give rise to most employment discrimination lawsuits," *id.* at 1068, misses the mark because much of the employment discrimination *problem* involves hiring discrimination. That few hiring suits are filed says more about the barriers to success in hiring-stage litigation than it does about the scope of the discrimination problem. *See* John J. Donohue III & Peter Siegelman, *The Changing Nature of Employment Discrimination Litigation*, 43 STAN. L. REV. 983, 1015–17 (1991).

⁶² Mitchell & Tetlock, *supra* note 13, at 1085 & n.204.

⁶³ For an excellent argument for the wrongfulness of what the author explicitly calls "Bayesian" discrimination, see Jody D. Armour, *Race Ipsa Loquitur: Of Reasonable Racists, Intelligent Bayesians, and Involuntary Negrophobes*, 46 STAN. L. REV. 781 (1994).

⁶⁴ *See* Alan David Freeman, *Legitimizing Racial Discrimination Through Antidiscrimination Law: A Critical Review of Supreme Court Doctrine*, 62 MINN. L. REV. 1049, 1052–54 (1978).

groups, Mitchell and Tetlock ask “what moral stance should be taken toward these individuals?”⁶⁵ When discussing measures of “modern racism,” they scornfully ask “whether people should be labeled racists for believing that many or most obstacles to racial equality are now internal to the Black community.”⁶⁶ And they show an intense concern for those accused of implicit bias by noting the “dangers of making false accusations of prejudice”⁶⁷ and by worrying that the IAT labels individuals with explicit egalitarian views as “implicit bigots on par with children reared in prejudiced households and taught to hold mean-spirited beliefs about minorities and to act out these prejudices.”⁶⁸

As Mitchell and Tetlock themselves acknowledge, these concerns are “ultimately a matter of political values, not scientific fact.”⁶⁹ If one sees antidiscrimination law as limited to punishing individuals who engage in irrational, animus-based discriminatory conduct, then Mitchell and Tetlock’s concerns about false accusations of prejudice will resonate. But advocates of using the law to counter implicit bias deny that antidiscrimination law should be about punishment or that it should focus only on irrational animus. Rather, they contend that antidiscrimination law should respond to the harmful effects of the “mechanisms of bias as produced by the current, *ordinary* workings of human brains—the mental states they create, the schemas they hold, and the behaviors they produce.”⁷⁰ They thus treat implicit bias not as a matter of individual fault, but as a social problem with harmful consequences to which the government must respond. The disagreement is, again, a normative one. It is not a scientific one.

In their conclusion, Mitchell and Tetlock pull back the curtain to reveal the normative agenda their “scientific” arguments really support. They argue that the claims of implicit bias advocates profoundly “shake the ontological foundations of American political culture” because those claims “imply that we cannot achieve equality of opportunity unless we have already achieved equality of result.”⁷¹ And they argue that those claims

⁶⁵ Mitchell & Tetlock, *supra* note 13, at 1083.

⁶⁶ *Id.* at 1082 n.190; *see also id.* at 1064 (arguing that explicit measures of “modern, symbolic, and aversive racism often include items that could easily serve as measures of ideological conservatism, traditional values, and the Protestant work ethic”). Note that this criticism reaches beyond implicit bias research to indict much research into *explicit* and *overt* prejudice.

⁶⁷ *Id.* at 1118; *see also id.* at 1101 (worrying about “false-accusation rates”).

⁶⁸ *Id.* at 1088; *see also id.* at 1104 (expressing concern that the IAT research program “is being used to indict the vast majority of the American population as implicitly prejudiced and disposed to discriminate whenever a suitable pretext emerges”).

⁶⁹ *Id.* at 1101.

⁷⁰ Kang & Banaji, *supra* note 7, at 1075 (emphasis added). Indeed, Kang and Banaji specifically disclaim the argument “that implicit bias-induced discrimination should produce the same legal liability as explicit animus-driven discrimination under current equal protection doctrine or federal antidiscrimination statutes.” *Id.* at 1077.

⁷¹ Mitchell & Tetlock, *supra* note 13, at 1120.

“stimulat[e] excessive suspicion of Whites among Blacks, suspicion that can crystallize into conspiracy theories that poison race relations.”⁷²

These are, of course, normative arguments. And there is substantial reason to doubt them. Terms like “equality of opportunity” and “equality of result” have strong political resonance, but they have no clear analytic meaning.⁷³ To the extent there is a meaningful distinction between the two concepts,⁷⁴ the claims of implicit bias scholars seem clearly to fall on the “equality of opportunity” side of the line. Those scholars do not claim that every race has to be equally represented in all professions; they claim that implicit and unconscious biases deny members of minority groups the opportunities to which they would be entitled if the biases were not present. To say that those scholars seek more than equal opportunity requires defining the term so narrowly as to guarantee nothing more than freedom from overt prejudice. If one disagrees that equal opportunity is so limited a concept, the bulk of Mitchell and Tetlock’s “scientific” arguments are basically irrelevant to an evaluation of the project of using the law to combat implicit bias.

III. THE NORMATIVE CHALLENGES THAT REMAIN

In the previous Part, I showed that Mitchell and Tetlock’s challenges are not so much scientific as normative. If one takes the view that antidiscrimination law should aim solely at punishing individuals who act on irrational racist and sexist animus, their arguments will resonate strongly. But if one takes the view that discrimination is a social problem that entrenches the subordinated and disadvantaged status of particular, socially salient groups, then their arguments will be far less relevant. In some respects, those arguments call for caution and further study before taking the more extreme steps urged by some scholars who write about implicit bias—steps like using the IAT as an employment screening device⁷⁵—but they do not undermine the strong reasons to believe that implicit bias contributes to a serious social problem that denies opportunities to blacks and other minorities in America. Indeed, by emphasizing the context-dependence of implicit bias, their arguments actually lend support to some of the proposals offered by legal scholars who urge a response to implicit bias.⁷⁶

I might just stop there. Mitchell and Tetlock have attempted to deliver a fatal blow to the implicit bias law-reform program, but their blow

⁷² *Id.* at 1119.

⁷³ See generally David A. Strauss, *The Illusory Distinction Between Equality of Opportunity and Equality of Result*, 34 WM. & MARY L. REV. 171 (1992).

⁷⁴ See *id.* at 172 (“The distinction between opportunity and result is an unhelpful and misleading way to categorize social institutions.”).

⁷⁵ See *supra* note 23 and accompanying text.

⁷⁶ See *supra* text accompanying note 61.

has largely missed the mark. Let us proceed with the program—with caution, to be sure, but let us proceed.

But things are not so simple. Mitchell and Tetlock's argument likely *will* resonate with many lawyers, judges, and policymakers. The justification that best fits antidiscrimination doctrine is one that focuses, not on the individual fault of the discriminator, but on the social harms of discrimination.⁷⁷ But most people—even most policy actors—have not internalized that justification. Legal interventions that go past the individual-fault justification and press against the social-harm one therefore encounter political and judicial backlash.⁷⁸ The adverse reaction to the Americans with Disabilities Act, whose requirement of “reasonable accommodation” centrally targets *rational* employer conduct, demonstrates this point.⁷⁹

What I have done in this Essay is therefore only the first step. Science does not defeat the implicit bias law-reform program, but science does not establish the case for that program, either. That program depends on a normative judgment that discrimination is not about fault but about a social problem—a normative judgment that is deeply contested among judges and policymakers. Unfortunately, arguments about implicit bias too often place heavy rhetorical emphasis on appeals to “science,” perhaps at the expense of developing the normative arguments for applying antidiscrimination law to that important phenomenon.⁸⁰

To be sure, the results of psychological studies can help to address the normative debate. The “shooter bias” studies, for example, are powerful because of the severe nature of the harm they suggest that implicit bias can cause.⁸¹ It hardly stretches widely accepted normative understandings to suggest that one's race should not be the basis for being shot—even if it is automatic processes of the brain, and not anything that is the “fault” of the officer, that causes a police officer to shoot an unarmed black man.⁸² When implicit bias affects actions that are not so immediately and

⁷⁷ See generally Bagenstos, *supra* note 21.

⁷⁸ See Bagenstos, *supra* note 1, at 40–45.

⁷⁹ See Linda Hamilton Krieger, *Socio-Legal Backlash*, 21 BERKELEY J. EMP. & LAB. L. 476 (2000).

⁸⁰ Kang and Banaji repeatedly seek to invoke the mantle of “science” in a way that, at least rhetorically, seems intended to do normative work, *see supra* note 17, although they also recognize that there are possible normative differences between discrimination actuated by implicit bias and other sorts of discrimination. *See supra* note 70. Linda Krieger recognizes that “even the best insights from the empirical social sciences cannot supply the normative principles needed for substantive lawmaking or resolve the conflicts between competing norms and interests so often implicated in the legislative and judicial processes.” Krieger & Fiske, *supra* note 6, at 1007. But her recent work on the issue rhetorically relies heavily on invocations of empiricism and science. *See generally* Greenwald & Krieger, *supra* note 3, *passim*; Krieger & Fiske, *supra* note 6, *passim*. Jolls and Sunstein touch on some of the normative concerns here, but only briefly. *See* Jolls & Sunstein, *supra* note 7, at 992–95.

⁸¹ *See supra* note 34.

⁸² *See* Banks et al., *supra* note 20, at 1173–75 (discussing how the shooter bias studies

ultimately consequential, but are instead “the everyday actions of virtually all of us,”⁸³ the normative case for responding to that bias becomes more complex.

So far, the most significant effort to develop a normative response to problems of implicit bias has come from those scholars who advocate a structural approach to antidiscrimination law.⁸⁴ Invoking the “democratic experimentalist” theory that merges early twentieth-century American pragmatism with late twentieth-century Japanese management principles, Susan Sturm has argued that the normative principles for responding to structural discrimination problems will best develop through a process in which those problems are addressed at the local level and information is shared through professional and other networks.⁸⁵ I gave qualified (though decidedly pessimistic) endorsement to that proposal in my earlier piece on these issues.⁸⁶

But to a large extent, the democratic experimentalist move is a punt. It remains incumbent on employment discrimination scholars to articulate and elaborate the normative theory under which the law should respond to implicit bias. Such theorizing can, of course, inform the efforts of localized activists who operate in a democratic-experimentalist world. But the theorizing is important, more fundamentally, to convince policymakers and judges to endorse the legal regimes that can address implicit bias.

My call, then, is for a renewed attention to antidiscrimination *theory*. As the social and behavioral sciences have gained prominence in the legal academy, much antidiscrimination work has come to explore the economics and psychology of discrimination. Work on implicit bias is just a part of that trend. There is much that this kind of work can do to help illuminate the important issues of antidiscrimination law. Just as did economic analysis, psychological research into implicit bias can help us better understand both the need for and the limits of antidiscrimination law. If continued implicit bias research shows increasingly strong correlations between measures of that bias and behaviors that (cumulatively) have significant discriminatory effects, that information will be centrally important to a full understanding of the discrimination problem. But at the same time scholars continue to develop that psychological research, we must continue to develop our normative understandings of antidiscrimination law to fit the findings of that research as well.⁸⁷ And that, ultimately, is the

were inspired by the widely publicized shooting of Amadou Diallo by New York police officers).

⁸³ Bagenstos, *supra* note 1, at 42–43.

⁸⁴ See Sturm, *supra* note 8, at 473–74, 559–64.

⁸⁵ See *id.* at 559–64.

⁸⁶ See Bagenstos, *supra* note 1, at 45–47.

⁸⁷ That, in part, was what I was trying to do in Bagenstos, *supra* note 21. Tristin Green has also taken up this challenge. See Tristin K. Green, *A Structural Approach as Antidiscrimination Mandate: Locating Employer Wrong*, 60 VAND. L. REV. (forthcoming 2007), available at http://papers.ssrn.com/sol3/papers.cfm?abstract_id=903791.

lesson of Mitchell and Tetlock's ostensibly scientific challenge to implicit bias research.

IV. CONCLUSION

Implicit bias research has already informed antidiscrimination scholarship, and it is likely to inform antidiscrimination scholarship even more in the future. But, as Mitchell and Tetlock's argument highlights, the implicit bias law-reform project is incomplete. This is true not so much for scientific reasons (though more work clearly needs to be done) but for normative ones. Mitchell and Tetlock's arguments underscore that point. Those arguments fail as a scientific challenge to implicit bias research because they largely rest on a set of unexpressed normative disagreements with implicit bias scholars. The implicit bias law-reform program still stands, but scholars advocating for that program must do a better job of articulating and defending the normative propositions on which their arguments depend. "Science" will not save antidiscrimination law. To save antidiscrimination law requires articulating and defending the normative principles that justify and guide the application of that law to newly understood forms of bias, whatever they are, and however they are discovered.

