

Can Restorative Justice Conferencing Reduce Recidivism? Evidence From the Make-it-Right Program*

Yotam Shem-Tov

Steven Raphael

Alissa Skog

September 8, 2021

Abstract

This paper studies the effect of a restorative justice intervention targeted at youth ages 13 to 17 facing felony charges of medium severity (e.g., burglary, assault). Eligible youths were randomly assigned to participate in the Make-it-Right (MIR) restorative justice program or to a control group in which they faced criminal prosecution. We estimate the effects of MIR on the likelihood that a youth will be rearrested in the four years following randomization. Assignment to MIR reduces the likelihood of a rearrest within six months by 19 percentage points, a 44 percent reduction relative to the control group. Moreover, the reduction in recidivism persists even four years after randomization. Thus, our estimates show that juvenile restorative justice conferencing can reduce recidivism among youth charged with relatively serious offenses and can be an effective alternative to traditional criminal justice practices.

*Yotam Shem-Tov, Assistant Professor, Department of Economics, University of California, Los Angeles; shem-tov@econ.ucla.edu. Steven Raphael, Professor and James D. Marver Chair in Public Policy, Goldman School of Public Policy University of California, Berkeley; and NBER; stevenraphael@berkeley.edu. Alissa Skog, Senior Research Associate, California Policy Lab, University of California, Berkeley; alissaskog@berkeley.edu. A pre-analysis plan for this project can be found at the Open Science Foundation, <https://osf.io/3sb4u/>. We are especially grateful to Katherine Miller, Maria McKee, Tara Regan Anderson, Mikaela Rabinowitz, and Jean Roland for the many insights they provided over the course of this project, as well as to the former District Attorney of San Francisco, George Gascon, who initiated Make-it-Right and advocated for the random allocation of individuals into the program which allowed for this study to be conducted. We are also thankful for invaluable information Lauren Brown and staff from Community Works West and Huckleberry Youth shared about the principles of restorative justice and the program's implementation. The California Police Lab provided continuous support and help throughout this project, we would especially like to thank April Chang, Johanna Lacoe, and Evan White. We thank Hadar Avivi, Annabelle Berrios, Gale Burford, Peng Ding, Sara Heller, Dmitri Koustas, Nicholas Li, Juliana Londoño-Vélez, Maxim Massenkoff, Conrad Miller, Andrew Penner, and Heather Strang for helpful comments and discussions. We also thank seminar and conference participants at Harvard Kennedy School Social Inequality Seminar, University of British Columbia, WEAI, UCLA, NBER SI Crime, and California Policy Lab for numerous helpful comments and suggestions. Chandni Raja, Shivani Ghatem, Logan Spencer provided outstanding research assistance. Any opinions and conclusions expressed in this paper are those of the authors and do not necessarily represent the views of the San Francisco District Attorney, San Francisco Juvenile Probation Department, Community Works West, Huckleberry Youth, or the California Policy Lab. We thank Arnold Ventures, the University of California Office of the President Multicampus Research Programs and Initiatives, MRP-19-600774 and M21PR3278, The James Irvine Foundation, and the Bylo Chacon Foundation for their generous support. The views expressed here do not necessarily reflect the views of the funders. All errors should be attributed to the authors.

1 Introduction

The United States criminal justice policy has historically relied on sanctions to enforce compliance. While the objectives of criminal justice policy are complex and multifaceted, much policy debate and research focuses on the efficacy of sanctions (e.g., [Kuziemko, 2013](#); [Aizer and Doyle, 2015](#); [Bhuller et al., 2020](#); [Rose and Shem-Tov, Forthcoming](#)). Sanctions became increasingly punitive beginning in the late 1970s, leading to large increases in correctional populations ([Raphael and Stoll, 2013](#); [National Research Council, 2014](#)). These policies had disparate impacts and generated large racial disparities in incarceration rates ([Neal and Rick, 2014](#); [Lofstrom and Raphael, 2016](#)). Recent years have seen an effort to dial back the severity of punishments and use alternatives to prosecution to reduce correctional population without impacting public safety ([Mueller-Smith and Schnepel, 2020](#); [Rose, 2020](#)). Alternative programs for juveniles might be especially effective, given the strong correlation between age and criminal involvement ([Hirschi and Gottfredson, 1983](#); [Ulmer and Steffensmeier, 2014](#)).

Restorative justice conferencing is an alternative that emphasizes accountability through repairing harm rather than imposing sanctions. While not always used instead of traditional prosecution, several countries and jurisdictions throughout the U.S. have experimented with this alternative prosecution path. Restorative justice programs typically involve a structured conference involving the victim, the accused, and supporters of both parties, followed by an agreement whereby the accused makes amends for the crime through a mutually agreed-upon set of actions. The process involves the individual charged with the offense taking responsibility for their actions and engaging in dialogue with the victim(s), as well as family and other community members, about the impact of their actions. The current evidence on the effectiveness of restorative justice programs in reducing recidivism is mixed ([Wilson et al., 2018](#)). Research evaluating such programs in the U.S. has relied mainly on observational methods. However, a key challenge is that participation in these programs is not random, and correlated unobservables are a likely threat to the identification of causal relationships.

This paper studies “Make-it-Right” (MIR), a restorative justice conferencing program implemented in San Francisco for teenagers charged with felony offenses of medium severity (e.g., burglary, theft, assault). Eligible cases were randomly assigned to either the MIR program (treatment group) or a regular felony prosecution (control group). Youth assigned to treatment were offered the opportunity to participate in the restorative justice program run by the organization Community Works West (CWW) instead of standard case processing and prosecution. Successful completion of the program results in the charges being permanently dismissed.

The program’s target population is high risk youth: 43 percent of control group members are rearrested within six months and 75 percent are rearrested within four years of treatment assignment. We estimate the intent-to-treat (ITT) and treatment-on-the-treated (TOT) effects of the program on recidivism in the four years subsequent to randomization. We find that MIR has

large crime-reducing effects. Juveniles assigned to MIR are less likely to be subsequently arrested for any reason, for felonies, or for more severe offenses. Correspondingly, after treatment, juveniles are also less likely to be convicted in the future. The effects of the program are large. Youths assigned to MIR were 19 percentage points less likely to be arrested for another offense in the six months after randomization. This estimate represents a 44 percent reduction relative to the rearrest rate among the control group of 43 percent. Moreover, the reduction persists even years after randomization. Those assigned to MIR are 15 percentage points (20 percent) less likely to be rearrested within three years and 27 percentage points less likely (32 percent) after four years. Among those assigned to MIR, 81 percent enrolled in the program, and 53 percent completed it. Accounting for imperfect take-up of treatment magnifies the results. Two-stage least squares (2SLS) estimates of the TOT with respect to enrollment and completion are roughly 1.3 and 1.9 times larger than the ITT effects.

Next, we compare the recidivism rates among the MIR control group and the full population of juveniles arrested for felony offenses in San Francisco between 2010 to 2020. Although eligibility to MIR was not random (e.g., there are restrictions based on criminal history, gang affiliation, and type of offense), the MIR control group’s rearrest rate is similar to that of the entire population of juveniles with felony arrests in San Francisco. Adjusting for differences in observable covariates using propensity score reweighting, the control group’s rearrest rate is slightly higher. Thus, the MIR experimental population is not a selected low-risk population. The similarity in risk between the experimental population and the broad population of juveniles arrested for a felony offense suggests that our experimental estimates may be predictive of the effect of expanding the program to the wider pool of juvenile cases.

Finally, we investigate the mechanisms through which assignment to MIR impacts recidivism. MIR has two essential components: the restorative justice conferencing itself and diversion from felony prosecution. As a diversion program MIR impacts case outcomes. For example, among youth assigned to MIR only five percent of the cases result in a felony conviction relative to 20 percent in the control group. To distinguish between MIR’s effects through restorative justice conferencing from any effects that are mediated through impacts on avoiding a felony conviction, we restrict attention to youth who have not been convicted of a felony offense due to the charges against them. Individuals whose case results in a felony conviction are likely those of equal or higher risk of recidivism. Among the control group, 20 percent of cases will result in a felony conviction relative to five percent in the treatment group. Thus, removing these cases drops more individuals of potentially higher recidivism risk from the control group. Nevertheless, in this restricted sample, the impacts of MIR on recidivism are similar to those in the full and non-restricted experimental sample. The relatively muted impact of the diversion component of MIR suggests that fostering empathy and dialogue with the victim are a driving force for MIR’s large recidivism-reducing effects. Moreover, it supports the hypothesis that the interaction with the victim during the restorative justice conference is a transformative experience for the youth.

Our results contribute to a broad literature across the social sciences on the determinants of criminal behavior. Previous studies document that individuals’ criminal behavior is influenced by punishment severity and swiftness (Drago et al., 2009; Hjalmarsson, 2009; Abrams, 2012; Aizer and Doyle, 2015; Davidson et al., 2019; Eren and Mocan, 2021), the likelihood of being caught (Di Tella and Schargrodsky, 2004; Blanes i Vidal and Mastrobuoni, 2018), the economic rewards of the crime (Draca et al., 2019; Kirchmaier et al., 2020), neighborhoods and peers (Ludwig et al., 2001; Kling et al., 2005; Damm and Dustmann, 2014), and opportunities in the labor market (Bushway, 2004; Raphael, 2014; Yang, 2017; Agan and Starr, 2018; Schnepel, 2018; Britto et al., 2020; Doleac and Hansen, 2020; Rose, 2021; Khanna et al., 2021).¹ However, little is known about whether or not criminal behavior can be influenced by less punitive interventions that use self reflection and empathy to directly change a youth’s decision-making.

Two related studies investigate the effects of cognitive-behavioral therapy (CBT) on anti-social and criminal behavior. Heller et al. (2017) conduct a randomized controlled trial (RCT) among youth in Chicago and show that CBT changes a youth’s decision-making, reduces criminal behavior, and increases high school graduation. Blattman et al. (2017) conduct an RCT in Liberia recruiting high-risk men aged 18-35. They find that CBT reduces anti-social behavior. However, the effects dissipate and are not significant after a year unless combined with a monetary grant. Our analysis complements these findings by showing that similar to CBT, interventions that appeal to one’s sense of responsibility (and perhaps shame) and, unlike CBT, specifically encourage pro-social behaviors have the potential to reduce recidivism. Our findings are especially encouraging when considering that our sample is comprised of high-risk juveniles charged with felony offenses of mid-level severity (e.g., burglary, assault) and not those charged with only minor offenses. Through the lens of Becker (1968), both MIR and CBT-type interventions do not change the economic incentives individuals are facing (i.e., the costs and benefits from criminal behavior); rather, they influence the individual’s decision making (i.e., preferences) while holding the setting fixed.

We also contribute to the growing literature that evaluates the efficacy of different diversion programs aimed at addressing the needs of those who become involved with the criminal justice system. Such diversion programs range from courts focused specifically on the treatment of defendants with severe mental illness (Cuellar et al., 2006), to drug courts that focus on compliance with substance use disorder treatment, sobriety, and recidivism (Mitchell et al., 2012), to specialty court focusing on the needs of specific populations such as veterans or domestic violence cases (Owens et al., 2021). Two recent studies use different sources of exogenous variation to identify the effects of diversion on recidivism. Augustine et al. (2021) exploit quasi-random assignment to judges and Mueller-Smith and Schnepel (2020) leverage natural experiments associated with shifts in diversion policy. Both studies find evidence that diversion programs reduce recidivism.

¹A growing recent literature also finds that parental criminal involvement can impact children’s likelihood of interacting with the criminal justice system (e.g., Bhuller et al., 2018; Dobbie et al., 2018; Huttunen et al., 2019; Arteaga, 2020; Norris et al., 2020).

Although related, restorative justice conferencing is fundamentally different than common diversion programs as it presents an alternative approach to addressing the harms caused by a criminal incident. For instance, the conference gives the victim an active and prominent role. Moreover, in our setting, we show that MIR’s large recidivism-reducing impacts are not mediated by its effects on avoiding a felony conviction (through its diversion component). Thus, our findings suggest that fostering dialogue (and perhaps empathy) between the victim and the accused can have meaningful long-term effects.

The remainder of this paper is organized as follows. Section 2 describes more in detail restorative justice practices and reviews the previous literature on such interventions. Section 3 describes the institutional setting, data sources, and summary statistics. Section 4 describes our empirical strategy. Section 5 presents our main findings and additional analyses regarding the external validity of the estimated effects of MIR. Section 6 discusses the potential mechanisms driving our results. Finally, Section 7 concludes.

2 Restorative Justice Conferencing: Theoretical Background and Empirical Evidence

In this section, we begin by briefly describing the theoretical considerations behind the current sentencing system in the U.S. and some of its drawbacks. We then explain what restorative justice is and its theoretical motivation. Lastly, we discuss the existing empirical evidence on the efficacy of restorative justice interventions in reducing recidivism.

2.1 The Current System of Punishment

Criminal justice policy and practice in the U.S. has historically been structured around several normative criteria. Consequentialist policy focuses on creating an enforcement and punishment architecture in service of deterrence, incapacitation, and rehabilitation. Policy levers devoted to increasing the likelihood of detection and the swiftness and/or severity of sanction are used to minimize criminal offending, often with an eye to the fiscal and social costs of exercising these policy levers. On the other hand, retributivist sentencing focuses on the moral content of the criminal offense, the degree to which the offense violates community norms, and the punishment proportional to the degree of harm caused by the offense, broadly defined. Retributivist theory focuses on the actions and autonomy of the individual, and the reaffirmation of community values by delivering through sentencing the “just deserts” merited by the violation of community norms. These alternative normative approaches to criminal sanctioning are in some instances complementary, but are often in conflict. For example, lengthy prisons sentences that incarcerate people into advanced ages who commit serious violent offenses in their youth may be justified based on the

gravity of the crime, but may not be justifiable in terms of deterrence and incapacitation (especially during the latter years of a sentence).

Despite their different objectives and at times prescriptions, these two normative approaches to criminal justice are firmly nested within policy approaches to criminal offending based almost entirely on sanctioning. Case adjudication centers around the actions and rights of the accused. Accountability, pursued in the name of either retributivist or consequentialist objectives, is enforced primarily through the severity of the sanctions issued.

Standard criminal case processing in the U.S. rarely addresses the underlying factors that contribute to one's propensity to criminally offend, or more generally, to have frequent interactions with law enforcement. For example, there is ample evidence of substantial overlap between individuals charged with and convicted of criminal offenses and individuals who have been victimized in the past ([Sampson and Lauritsen, 1990](#); [Berg et al., 2012](#)). Moreover, the criminal justice system often fails to center the experiences and needs of crime victims. Criminal offenses are framed as crimes against the community or state, with the actual victim serving primarily as a witness. While convicted offenders can be ordered to pay some form of restitution and victims may seek redress through civil litigation beyond the criminal case, standard processing has very little restorative content from the victim's viewpoint beyond metaphorically making the accused pay for the harm they have caused.

2.2 Restorative Justice

Advocates of restorative justice argue for the direct involvement of the harmed and a fuller consideration of the factors contributing to the behavior of the person committing the harm. Such advocates argue that criminal cases can be adjudicated in a manner that builds greater empathy for the victim and appreciation of the consequences of one's actions while simultaneously improving the well-being and satisfaction of victims. As summarized by [Umbreit and Armour \(2011\)](#), adversarial criminal case processing focuses on establishing what laws have been broken, who broke them, and what punishment the person(s) responsible deserve(s). By contrast, restorative justice systems focus on establishing who has been harmed, what needs to occur to restore the welfare of the harmed, and who is obligated to make the restoration.

While restorative justice programming can take many forms, in criminal justice settings it usually involves direct conferencing between the person responsible for the criminal offense and the person who has been harmed.² Restorative justice conferences can occur in addition to standard criminal processing ([Angel et al., 2014](#)) or in lieu of criminal processing (as in MIR). The process

²[Sherman and Strang \(2007\)](#) discuss several alternative incarnations of restorative practices including direct conferencing, conferencing through third-party mediation, and the payment of restitution. The authors note the long history of restorative practices in response to criminal offense throughout the world, and the influence of the practice of indigenous cultures (e.g., the Maori in New Zealand, Native Americans in the United States) in informing conventional thinking about restorative justice.

may involve a direct discussion between victim and perpetrator of a specific offense (McGarrell, 2001) or the conference may be between individuals who have committed and have been harmed by specific offenses but not involved in the same incident (Fulkerson, 2001). Conferencing requires the consent of both the harmed as well as the person who committed the offense, and typically requires that the offender acknowledge their culpability. To date, much experimentation has occurred in juvenile justice settings. These settings are seen as early opportunities for alternative case processing that may prevent or reduce adult criminal offending.³ Moreover, while these methods have a longer history of application in juvenile settings and for less serious offenses, some countries have experimented with application to more serious criminal cases, including domestic violence (Strang et al., 2002).

A restorative justice conference usually begins with a statement by the person accused of the crime, followed by an opportunity for the harmed party to directly address the accused. Conferences also typically incorporate supporters of both the harmed and the accused, with everyone involved permitted to communicate openly about how the offense impacted their well-being. Supporters often include family members and friends; in juvenile justice settings the family support system almost always includes a parent or guardian. Conferences close with an agreement phase, whereby the victim and the person who committed the offense agree to a plan of action that the person will undertake to restore the welfare of the harmed. The actions can include writing formal letters of apology, paying restitution, agreeing to specific community service, and/or tailored actions of good faith.⁴

Restorative justice has its origins in a theory of shaming in criminal justice practice first articulated by Braithwaite (1989). Braithwaite posits two forms of shaming for those who commit harm. Stigmatic shaming permanently associates the criminal offense with the character of the individual who committed the offense and thus creates separation between the individual and society. By contrast, reintegrative shaming separates the offense from the individual, effectively condemning the action but offering the person a path back into the community through the acknowledgement of responsibility and efforts to make amends. While modern incarnations of restorative justice tend to emphasize restoration over shaming, most do require that the person who committed the offense take responsibility for their actions. Moreover, through direct communication during the conference, the person who committed the offense is typically confronted with the consequences of their actions not only for the victim, but often by the victim's family and friends, and even members of their own family. In theory, this process should not only engender greater empathy for the person harmed, but will also make the accused more mindful and aware of how their actions

³The U.K. has also experimented with adjunct restorative justice conferences in adult corrections (Shapland et al., 2008).

⁴The authors observed a presentation of the service provider for MIR, CWW, to the San Francisco Sentencing Commission that highlighted a case where it was revealed during the conference that the youth who committed an offense was a talented muralist and the victim a lover of Disney characters. An element of the action plan from the conference was that the youth would paint a mural of the harmed party's favorite character in her apartment, a provision of the agreement on which the youth willfully and enthusiastically followed through.

impact others. Beyond the impact on the person causing harm, many hypothesize that the ability to conference directly with the person causing harm and the efforts towards sincere restoration minimize the trauma suffered by victims and permit them to move on with their lives to a greater degree than would otherwise be possible.

2.3 The Empirical Evidence on the Efficacy of Restorative Justice Conferencing

While there have been only a handful of small-scale restorative justice programs in the U.S. juvenile justice system, restorative justice is a key component of juvenile justice in New Zealand ([Ministry of Justice, 2004](#)) as well as a standard juvenile diversion alternative used in Australia ([Little et al., 2018](#)) following a decade of experimentation ([Strang et al., 2013](#)). The U.K. has also experimented with restorative justice interventions, usually in addition to standard criminal case processing rather than as an alternative ([Shapland et al., 2008](#)).

The existing evaluation research focuses both on victim outcomes (usually measures of satisfaction, or post-incident indicators of trauma) as well as measures of recidivism among those who committed the offenses. The extant body of research generally finds that restorative justice programs are well received by victims and may help alleviate post-traumatic stress symptoms ([McCold and Wachtel, 1998](#); [McGarrell, 2001](#); [Angel et al., 2014](#); [Brooks, 2013](#); [Sherman et al., 2015](#)).⁵

In contrast to findings regarding the effects on victims, evidence on the impacts of restorative justice interventions on recidivism is less conclusive. In a relatively recent review, [Livingstone et al. \(2013\)](#) highlights there is “currently a lack of high-quality evidence regarding the effectiveness of restorative justice conferencing for young offenders” and more research is needed. The Australian experiments reviewed by [Sherman et al. \(2015\)](#) find some evidence of a reduction in repeat offending, especially for offenses that involved a victim. However, there is some weak evidence in sub-group analyses of increased offending, though the estimates are underpowered and appear to be exploratory. An important feature of these experiments is that randomization occurred only among offenders who agreed to participate, leading to a selected sample. The U.K. also experimented with restorative justice, usually as a complement to criminal prosecution rather than as a substitute. In conjunction with the experimental evidence from the U.K., existing research seems to indicate larger crime-mitigating impacts when applied to individuals who have committed more

⁵[Fulkerson \(2001\)](#) studies an offshoot of restorative justice conferences on domestic violence victims. The author investigate the effect of victim impact panels (VIPs) in domestic violence cases in Arkansas. VIPs permit crime victims to directly address offenders, usually as a group, in a safe and secure and environment that is not restricted by the formal procedures of a courtroom. Treatment subjects were asked to participate in a VIPs whereby a panel of victims were given the opportunity to address a panel of offenders, with the one restriction being that victims could not participate on a panel where the individual that victimized them was among the group of offenders being addressed. While the study reports that most participants believed that the panels were helpful (both among victims and offenders), no contrast is drawn with individuals assigned to the control group.

severe offenses (Sherman and Strang, 2007). Moreover, Bonta et al. (2002) provide an observational evaluation of a restorative justice program that took place in the 1990s in Canada and acted as an alternative to incarceration for convicted offenders. They find large reductions in recidivism. However, the lack of random or quasi-random treatment assignment implies the findings could also be driven by selection and unobserved differences between the treated and control participants.

In the U.S., the evidence is mixed. First, in an observational study Brooks (2013) evaluates a restorative justice community conferencing program focused on juveniles in Baltimore, Maryland. Brooks uses a propensity score matching design to compare youth who participated in the program to observably similar individuals who did not. Interestingly, Brooks finds higher recidivism rates among the participants of the restorative justice program.

Second, the results of two RCTs from the 1990s are inconclusive and find no persistent reductions in recidivism. McGarrell (2001) and McGarrell and Hipple (2007) evaluate an RCT of a restorative justice program in Indiana focused on children (the average age is 12.5). They find meaningful decreases in recidivism within one year relative to a control group allocated to another diversion program. However, in follow-up work examining the long-run effects of the intervention, they find that the intervention affected neither the likelihood of a rearrest nor the time to a new arrest (Jeong et al., 2012). McCold and Wachtel (1998) evaluate a family-group conference intervention in Pennsylvania and find no evidence of reductions in recidivism after one year. However, the take-up among those assigned to treatment was relatively low (roughly 40 percent).

The two experiments studies by McGarrell (2001) and McCold and Wachtel (1998) are more related to our evaluation of MIR than the non-U.S. studies reviewed in Sherman et al. (2015). In these two RCTs, randomization was done without conditioning on the youth’s willingness to participate (similar to MIR and unlike the non-U.S. studies), and they involve only juveniles who have not been convicted. However, they also differ from MIR in several key aspects. In McCold and Wachtel (1998), the conferencing was conducted by *police officers* and involved only youth charged with infractions and misdemeanors (no felonies). In McGarrell (2001), youth were especially young and generally engaged in less severe offenses. Another key difference is that MIR has an agreement monitor that supports the youth in completing the agreement post-conferencing. The agreement monitor joins the restorative justice process at the conference and then meets with the youth regularly (at least once a week) post-conferencing to make sure they are on track to fulfill the agreement signed during the conference.

Our evaluation of the MIR program contributes to this literature in several ways. First, it provides credible evidence that restorative justice conferencing can cause lasting reductions in recidivism. Importantly, there are no observed or unobserved confounders to the intervention in our setting since assignment to the treatment and control groups was done at random. Moreover, the take-up was high; 81 percent of those assigned to MIR enrolled in the program, and 67 percent completed the program among those who enrolled. These rates are especially high given that the accused’s consent to participate in MIR was not an eligibility requirement for inclusion in the

randomization. Second, unlike other restorative justice experiments where youth assigned to the control group were funneled into various alternative diversion programs, control group members under MIR faced traditional felony prosecution. Third, it highlights the potential importance of having an agreement monitor to check-in and support the youth after the restorative justice conference is over, and the onus of completing the agreement is on the youth.

A key difference between MIR and most non-U.S. experiments is that the accused is not required to consent to participate in the RCT—that is, consent is not an eligibility requirement for inclusion in the randomization. Therefore, the ITT estimates take into account any non-participation due to a lack of willingness on the accused’s side. Moreover, effects for those who participated in MIR can also be recovered using 2SLS, as discussed below.⁶

Thus, our ITT estimates capture exactly the policy-relevant estimand of comparing assignment to a restorative justice intervention relative to a standard felony prosecution. Lastly, we find economically meaningful and statistically significant estimates that MIR reduces recidivism both in the short-run and in the long-run after four years from randomization. These findings are informative to the current debate on whether restorative justice conferencing can reduce recidivism or not.

3 Setting and Data

We begin by describing our data, the details of the MIR program, and the process by which juveniles were randomized to either the MIR treatment group or a business-as-usual control group. We also provide summary statistics on the MIR experimental sample and compare it to the full population of juveniles arrested for a felony offense in San Francisco.

3.1 Data Sources

Our evaluation draws upon three data sources, all pertaining to criminal offenses occurring within San Francisco. First, we were provided programmatic information on all youth who were part of the MIR experiment. Programmatic data include information on whether the eligible youth was assigned to MIR or the control group and information on MIR program participation. For those who enrolled in MIR, we were provided information on key dates including enrollment date, conference date, date of completion, and for those who did not complete the program, the date the case was sent back to the office of the San Francisco District Attorney (SFDA). Second, we were provided data on the universe of juvenile arrests for offenses occurring in San Francisco between October 2010 and November 2020. The data includes a description of the offense (by penal code), demographics such as sex, race, date-of-birth, name, and age at first criminal incident. These

⁶In MIR, the consent of the victim was a requirement for eligibility to be randomized, however, in conversations with the handling juvenile prosecutor we have been told that all the victims consented.

records allow us to construct measures of criminal history for the youth participating in the RCT. Third, we were provided with administrative records on all adults arrested and presented for a charging decision in San Francisco between October 2010 and November 2020. Given that some of the youth in the experiment turn 18 shortly after the arrest for the offense associated with their inclusion in the experiment, the adult arrests records are essential to construct accurate recidivism outcomes.

3.2 The Make-it-Right Program

The SFDA piloted the MIR program at the end of 2013. The program diverts youth who have been arrested for certain felony offenses to a restorative justice conferencing program. Conferencing involves facilitated, community-based conversation between the involved minor, their family, the person harmed, and a community representative, leading to an agreed-upon plan for addressing that harm. The SFDA permanently drops the criminal charges for youth who successfully complete MIR, and, therefore, these juveniles will not have a conviction for this case on their record. Importantly, MIR is a pre-charging diversion program. Eligible youth are randomized to either the treatment or control groups after the juvenile prosecutor reached the decision to file charges but before charges have been filed. Thus, all the individuals in the control group will face criminal charges and youth assigned to MIR but who did not complete the program will also automatically face criminal charges.

Eligibility criteria. The program is targeted at juveniles 13 to 17 years of age charged with medium-severity felony offenses such as the unlawful taking of a vehicle, grand theft, burglary, or assault. Thus, the sample of eligible youth includes only those arrested for felony offenses, though the program potentially makes exceptions for youth classified as repeat misdemeanor offenders. To be referred to the program, the youth must be a resident of San Francisco County or Northern Alameda County, must have no prior 707(b) arrests or sustained petitions.⁷ In addition, the youth must not have caused moderate or significant injuries to the victim, cannot be affiliated with a gang, cannot have used a weapon, and the minor cannot be currently under probation supervision or in detention.⁸

MIR was intended to be a relatively small pilot: SFDA expected to enroll no more than 25 individuals per year. In practice, the number of eligible youth and enrollees was lower than expected. One reason for this is the steady reduction in juvenile crime in San Francisco, which

⁷707(b) offenses are those that would count as a strike under California’s three-strikes for juveniles 16 and over, for example, violent felonies such as kidnapping or robbery.

⁸These requirements ensure that an eligible youth does not have any other pending cases; otherwise, they would be either on probation supervision or in detention. Moreover, if the youth is rearrested for a new offense, then either their participation in MIR will be stopped, or the new offense will be merged with the original one for which the youth was assigned to MIR. Thus, by construction, youth randomized into MIR will not be a part of any criminal justice proceedings during their participation in MIR.

affected the overall volume of cases. In total, the MIR RCT lasted 5.5 years and included 99 participants in the restorative justice treatment group and 44 participants in the control group at the time it concluded in May 2019.

Randomization Procedure. The SFDA’s Office designed and led the implementation of the MIR RCT; the research team was not involved in devising nor implementing/monitoring the randomization process. The SFDA used block randomization to create treatment and control groups, with individuals randomized at the case level, which corresponds to individuals except for cases involving co-defendants. The randomization process was designed to separate the actual assignment from the prosecutor in charge of the program. Specifically, once a case was deemed eligible, it was sent to a paralegal not involved with the program and unrelated to the juvenile prosecutor, who consulted a prepared list of assignments, selected the next available assignment and then communicated the assignment back to the juvenile prosecutor.⁹

The MIR program. Table 1 provides a concise description of the MIR program, its activities, and each of its steps. Once assigned to treatment, CWW, a non-governmental organization located in Oakland, California, specializing in restorative justice, assesses the youth’s ability to participate.¹⁰ If the youth is deemed unsuitable, the case is referred back to the SFDA for traditional prosecution. An essential requirement for participation is that the youth demonstrates reflection and accountability for his/her actions. Minors and their parents may decline to participate, effectively opting for the case to be referred back to the SFDA for prosecution. For cases that proceed to conferencing, CWW does all of the pre-conference planning involving the youth (referred to as the responsible party), the victim (harmed party), and any other individuals who will take part in the conferencing (such as parents and/or supporters of the harmed and responsible parties). Moreover, CWW also mediates the conference. Post-conference case management and compliance monitoring is managed by the Huckleberry Youth’s Community Assessment and Resource Center. Youth who fail to follow through with the program have their cases referred back to SFDA for prosecution.

Juveniles in the control group continue through the traditional juvenile justice process. They are charged and prosecuted in juvenile court, with the minor being supervised by the JPD during the trial. While some minors are detained while their case is processing, most are supervised in the community using electronic monitoring or day reporting centers.

⁹Initially, the SFDA randomly assigned 50 percent of eligible individuals to MIR. However, shortly after commencing the experiment, the SFDA discovered that the number of youth eligible for MIR was lower than initial estimates as overall juvenile crime and incarceration in the County was decreasing. The assignment probability to treatment was thus altered to 70 percent in May 2014 (six months after the pilot date). As we discuss below, including cohort fixed effects yields almost identical estimates.

¹⁰The SFDA sought consent from all victims before randomizing eligible youth. The assistant district attorney in charge of the program from its start informed us that all victims consented to have their case included in the study.

Figure 1 depicts the flow of cases through the two treatment arms. In total, 143 cases were deemed eligible and randomly assigned, with 99 (69.2 percent of study subjects) assigned to the treatment group and 44 assigned to the control group (30.8 percent). Youth assigned to MIR either enroll in the program or are deemed unsuitable (e.g., when refusing to assume responsibility for their share in the incident). The take-up rate was quite high, with 80.8 percent of those assigned to MIR enrolling in the program. This contrasts with the earlier restorative juvenile justice experiments in the U.S. that we review above. The relatively higher take-up rate may reflect the fact the MIR enrolls youth charged with relatively serious offenses (all felonies) relative to the requirements of lower-level charges for eligibility in prior experiments, and who face prosecution and potentially severe sanction as the alternative to participation in conferencing. The process of enrolling into MIR was relatively quick. On average, juveniles waited 15 days and the median waiting time was 21 days.

Among those enrolled in MIR, 66.7 percent completed the program and ultimately had their cases dropped. The average and the median duration of the program (time between enrollment and completion) were six months. Overall, 52.5 percent of those assigned to MIR completed the program. Thus, while most youth assigned to MIR completed the program, a non-negligible portion did not. Multiple reasons can lead to this result. For example, the accused and the victim did not reach an agreement during conferencing, or the youth did not fulfill the agreed contract with the victim.

In what follows, we present the final evaluation of the MIR program that makes use of the full sample of juveniles randomized to the program between 2014 and 2019.¹¹

3.3 Summary Statistics

Table 2 presents summary statistics for the different groupings of juveniles randomized under the MIR experiment and for all juveniles referred to JPD for felony offenses from October 2010 through November 2020. The first column presents average characteristics for youth assigned to the control group, while the second and third columns present summary statistics for all youth assigned to MIR, and youth who enrolled among those assigned to MIR (compliers), respectively. The final column presents descriptive statistics for all juveniles referred to JPD for felony offenses.

Comparisons of the averages for the control (Column (1)) and treatment (Column (2)) groups reveal that random assignment yielded balance on most observable characteristics. The square brackets in Column (2) report p-values for the null hypothesis that the averages in Columns (1) and (2) are equal. While there are a few differences, an F-test of the overall significance of a

¹¹In an early analysis of the program we study, [Huntington et al. \(2017\)](#) provide a preliminary outcome evaluation of the MIR program. This earlier intervention was severely under-powered, and focused largely on reviewing comparable experiments, thinking through the decision of whether to continue the experiment, and some very preliminary outcome comparisons. Our analysis was completed independently of this effort, based on separately procured data extracts, and a completely independent process of data pre-processing and analysis.

regression of a treatment dummy on the covariate list yields a p-value of 0.757. Moreover, all the individual p-values are not significant at the five percent level, and only two are significant at the ten percent level. Roughly 90 percent of youth in the study are male, approximately half are Black, and one-third are Hispanic. Average age at arrest is 16, though age at first arrest is lower for both groups (14.75 for the control group and 15.2 for the treatment group). Around one-fourth of the treatment and control groups have a prior arrest in San Francisco. However, prior arrests for felony offenses are less common (7 percent among the control group and 14 percent in the treatment group). Moreover, the average number of prior arrests is low among both the control and treatment groups, 0.7 and 0.37, respectively. Although there are some minor (and non-significant) differences in balance, adjusting for covariates does not change our estimates as we discuss in Section 5.3. Moreover, any imbalances are not systematic; for example, individuals in the control group are less likely to have a prior felony arrest; however, they have more prior arrests on average.

Next, we describe the type of offenses youth in MIR are typically facing. The most common charge is felony theft (64 percent of the control group and 66 percent of the treatment group), followed by burglary (approximately 40 percent) and felony assault (roughly 14 percent).¹²

The final four rows present means for a predicted recidivism (defined as a subsequent arrest in San Francisco) variable for the control and treatment groups. To generate these predicted values, we use auxiliary data on the entire sample of youth arrested for felony offenses between 2010 and 2020 not including youth who participated in the RCT. We estimate an OLS model where the dependent variable is an indicator for a rearrest for any new offense within t months of the controlling incident. The explanatory variables include all of the covariates presented in Table 2. We then use the estimated parameters from this model to generate predicted recidivism probabilities for each youth in the MIR treatment and control groups. Appendix B describes in more detail the construction of the predicted recidivism variables. Roughly 30 percent of the treatment and control groups are predicted to be arrested on a new offense within six months, with predicted felony recidivism rates of approximately 20 percent for both groups.

Comparison of the means in Columns (2) and (3) suggest that youth assigned to the MIR program and who also enrolled are indistinguishable from youth who do not take up (as can be seen from the p-value of the joint F-test presented at the bottom of Table 2). Moreover, as summary measures, the predicted recidivism probabilities both for any new offense and for felony offenses are similar for compliers and never-takers (those assigned to MIR who did not enroll into the program).¹³

Relative to the broader population of juveniles referred to JPD, described in Column (4), the

¹²Note, the charge proportion sum to greater than one since a given case is often associated with more than a single charge.

¹³In addition, the averages of 12-month any-recidivism and felony-recidivism predictions for the never-takers are 48.8 and 33.3 percent. These predicted recidivism values are slightly higher than the compliers means of 45.3 and 29.9 percent.

MIR youth are more likely to be male, somewhat less likely to be Black or Hispanic, are less likely to have a prior arrest, and have a first arrest occurring at a slightly older age (approximately 15 vs. 14.23). Moreover, there are notable differences in the charge distribution, reflecting the MIR eligibility criteria. Specifically, over a third of non-MIR juveniles are arrested for severe person offenses (1.8 percent for homicide, 1.3 percent for a sex offense, and 34.1 percent for robbery) compared with roughly two percent of MIR youth arrested for robbery. We also observe a higher proportion of all JPD referrals involving a weapons offense relative to MIR youth. Finally, the predicted recidivism measures are lower for the MIR youth.

To give more context on our setting and who are the individuals interacting in the restorative justice conferencing, Appendix Table A.1 presents the demographic characteristics of the victim (harmed party) and compares them to those of the accused (responsible party).¹⁴ A few patterns are worth noting. First, the average age of the victim is 35, roughly double the average age of the average accused youth. Second, the racial composition of the two groups is meaningfully different, and in 79 percent of the incidents, the harmed party is of another race/ethnicity than the responsible party. While most of the youth are Black and Hispanic, most of the individuals in the harmed group are White and Asian. Females are also overly represented in the harmed (victim) group, 41 relative to 11 percent among the responsible youth.

4 Empirical Strategy

Our empirical approach is straightforward given that individuals were randomized between assignment to MIR or the control group. Comparisons of mean outcomes among units assigned to treatment and control regimes are sufficient to identify the ITT effect of being assigned to MIR relative to felony prosecution (control regime). We begin with a set of empirical estimates that compare mean outcomes of individuals assigned to MIR (the treatment) relative to those assigned to control. We estimate the ITT effect using Equation (1) which describes the relationship between assignment to MIR and the likelihood of being arrested on a new charge within t months from the date of randomization:

$$Y_{it} = \gamma_0 + \gamma_1(\text{Assigned MIR})_i + \epsilon_{it} \quad (1)$$

where $Y_{it} \in \{0, 1\}$ indicates whether the youth is arrested on a new charge within t months. The individuals in our sample have been randomized to MIR or the control regime in different dates and hence we observe less individuals for longer time horizons. To exploit all the information in our data we complement Equation (1) with Kaplan-Meier estimates of the failure function (defining a failure as an arrest for a new charge).

¹⁴The information on the demographic characteristics of the victims in the MIR experimental sample (although partial) was provided to us by CWW and was collected as part of the restorative justice conferencing process.

In addition to the ITT analyses, we also estimate the TOT effect of participation in MIR, by using the random assignment to MIR as an instrumental variable for whether or not a youth enrolled into MIR. Equations (2) and (3) describe our 2SLS estimator:

$$Y_{it} = \beta_0 + \beta_1(\text{Enrolled MIR})_i + \eta_{it} \quad (2)$$

$$(\text{Enrolled MIR})_i = \alpha + \alpha_1(\text{Assigned MIR})_i + \xi_i \quad (3)$$

The 2SLS estimator from Equations (2) and (3) identifies the TOT under the LATE framework assumptions (Imbens and Angrist, 1994; Angrist et al., 1996), as there is only one-sided non-compliance in our setting (Bloom, 1984).¹⁵ To quantify the magnitude of the estimated ITT and TOT effects, we compare them to the control group complier mean (Katz et al., 2001).

Throughout our analysis, we report two types of p-values. The first is based on cluster-robust standard errors, clustered at the case level. The second p-value is based on randomization inferences using random permutations of cases to placebo MIR and control regimes to generate the sampling distribution under the null hypothesis (as was suggested by Young, 2019). This procedure is known as randomization, Fisherian, or permutation inference and was first proposed by Fisher et al. (1935). As was advocated by Chung and Romano (2013), we use as our test statistic the standardized t-statistic.¹⁶

Finally, in our pre-analysis plan we specified that we will evaluate only one-sided hypothesis tests pertaining to whether the MIR program reduced the likelihood of recidivism or not. This choice was aimed to maximize statistical power given our small sample size.¹⁷ Pre-analysis plans have been mentioned in the literature on research transparency and reproducibility as a tool to increase statistical power by pre-specifying one-sided hypothesis (e.g., Olken, 2015).¹⁸ Given our pre-specification choice of focusing on one-sided hypothesis tests, we are unable in this study to test whether the program increases the likelihood of an arrest on a new offense among participants.

¹⁵In our setting, these assumptions are likely to hold. Monotonicity must hold since all the individuals in the control group did not participate in MIR. Exclusion is also likely to hold as assignment to MIR is unlikely to impact recidivism except through participation in MIR.

¹⁶Note that randomization inference provides valid p-values for two null hypotheses. First, it provides finite-sample exact inference on the sharp null of no treatment effect on all individuals. Second, when the sharp null hypothesis is incorrect (e.g., some units are impacted positively while others negatively), randomization inference still provides asymptotically valid inference on the null hypothesis of no average treatment effect, also known as the Neyman null hypothesis. Thus, regardless of whether the sharp the null hypothesis is correct or not, randomization inference provides valid asymptotic inference on the null of no average treatment effect. For the latter interpretation to hold, a standardized statistic (e.g., t-statistic) must be used as the test statistic in the randomization inference procedure. The second interpretation of randomization inference tests was formalized only recently by Li and Ding (2017), Wu and Ding (2020), and Zhao and Ding (2021).

¹⁷We submitted the pre-analysis plan before looking at the data and without knowing the exact sample size. However, we knew its general range.

¹⁸Another example is Christensen and Miguel (2018) who advocated the use of pre-analysis plans and mentioned as one of their advantages the fact that they can allow researchers to specify their interest in one-sided hypothesis in advance and, by doing so, increase the accuracy of statistical tests, “PAP [Pre-Analysis-Plans] bind the hands of researchers and greatly limit specification searching, allowing them to take full advantage of the power of their statistical tests (even making one-sided tests reasonable).”

However, this is unlikely since the signs on all our estimated effects point towards reductions in recidivism. Although we find no evidence that MIR increases recidivism, this is not always the case. As we noted above, in some instances, prior experimental results from Australia (Sherman et al., 2015) revealed situations where restorative justice programs increased offending among some individuals.

5 Main Results

We begin with ITT estimates by visualizing the recidivism patterns that we observe in the data among individuals who have been randomly assigned to MIR and those assigned to the control regime. Figure 2 presents Kaplan-Meier estimates of the failure functions depicting the relationship between the probability of being rearrested on a new offense at least once and the number of days since randomization took place over a four-year period.¹⁹ For both groups, the likelihood of rearrests is quite high, with nearly half of the control group arrested within six months of randomization and over 70 percent rearrested by the end of the four-year period. However, we observe markedly lower rearrest rates among youth assigned to the treatment group. The difference in the percentage ever arrested on a new offense reaches roughly 20 percentage points within six months. It then fluctuates around this level for the remainder of the observation period (though we should note, the number of observable individuals for the latter periods narrows to roughly half of our sample).

Future arrests that result in a conviction are considerably less common for both the treatment and control groups. Nonetheless, a sizable part of the control group is subsequently convicted: one quarter is convicted for a new offense within six months of randomization, and over 40 percent is convicted by the end of the four-year period. Thus, again, we see lower failure rates for youth assigned to the treatment group, with about a ten percentage point gap opening up within six months and remaining at approximately this level over the four-year observation window. Finally, we examine the effects of MIR on measures of recidivism that focus only on more severe interactions. Appendix Figure A.1 shows similar results when measuring recidivism as future felony arrests or as future arrests for offenses that are at least as severe as the original charges for which the youth was arrested. These findings show that MIR causes reductions in recidivism across a variety of different measures, reducing both the likelihood of any future arrests as well as the likelihood of arrests for relatively serious offenses.

We perform formal hypothesis tests for equality of the two cumulative failure functions using

¹⁹Note, the last randomization occurred in October 2019, and we observe recidivism data through (including) November 2020. Hence, we have at least 14 months post-randomization for all youth, but obviously longer periods for youth randomized into the program in the early years. Figure 2 presents Kaplan-Meier estimates of the failure function which are based on one minus the product of period-specific survival probabilities (through a given time interval), implicitly assuming that the recidivism hazard function is stable across cohorts defined by randomization date.

the standard non-parametric Peto-Peto-Prentice test (see [Klein and Moeschberger \(2006\)](#) for a textbook description) and calculate p-values in two ways. First, using a non-parametric inference procedure that is based on asymptotic approximations (as is standard practice). Second, we use randomization inference and calculate the finite sample distribution of the test statistic from 1,000 random placebo permutations of the treatment assignment (also known as Fisherian or permutation inference, [Fisher et al. \(1935\)](#)). The top of Figure 2 reports the one-tailed p-value from these tests. For all arrests (Panel (a)), both tests reject the null hypothesis of equal failure functions, with the p-value from the randomization inference ($p = 0.017$) roughly 2.5 times the value based on the standard non-parametric asymptotic inference ($p = 0.0071$).²⁰ For arrests that lead to a conviction (Panel (b)), we observe marginal significance ($p = 0.072$) based on the non-parametric test and a p-value of 0.111 using randomization inference.

Table 3 presents our principal results. Each column presents model estimates where the dependent variable is an indicator variable equal to one if the youth is arrested on a new offense within the time period indicated by the column heading (e.g., within 6 months, within 12 months and so on). Panel (a) presents our 2SLS estimates of the TOT effect of participating in MIR on subsequent rearrests. Along with each estimate, we report standard errors clustered at the case level (in parentheses), the p-value from a one-tailed test of the null hypothesis of a non-negative impact on recidivism against the alternative of a negative (i.e., crime reducing) impact based on the standard error (in curly brackets), and the p-value from a comparable one-tailed test using 1,000 placebo permutations of the treatment assignment (in square brackets). Panel (b) presents the ITT effects along with the comparable standard error and p-value calculations to those presented for the TOT estimates. Panel (c) presents the first stage results for the effect of assignment to MIR on participation. Given that there is no crossing over from control to treatment in the data, these first stage effect estimates are the proportion of youth who are assigned to MIR and who comply and enroll in the program. Finally, at the bottom of the table we present averages of the dependent variable outcome for the control group and for compliers under the control regime.

Before describing the effect sizes, we should note the high rates of rearrest observed among the control group. 43 percent of control group members are rearrested within six months of random assignment, 63 percent are rearrested within two years, and 75 percent are rearrested within four years of assignment. The patterns in these average control outcomes highlight the fact that MIR is having an impact on new arrest rates for youth who are at a very high risk of future contact with the criminal justice system.

A second pattern to note is the high take-up rate. While the take-up rate varies slightly across the time periods analyzed, it never falls below 73 percent (Panel (c) of Table 3).²¹ Again, this

²⁰One reason for the fact that randomization inference yields larger p-values relative to the standard inference formulas is that the Peto-Peto-Prentice test p-value does not take into account clustering at the case level in the assignment to treatment. As we show in Table 3, once clustering is taken into account, randomization inference and regular cluster-robust standard errors yield similar p-values.

²¹Note, differences in sample size across models are the result of differences in the post-randomization time horizon

stands in contrast to the juvenile interventions reviewed above where take-up rates are considerably lower and we speculate that the relatively high take-up rate for MIR is likely driven by the fact that the alternative is a felony prosecution.

The reduced form effects (i.e., the ITT effect estimates) indicate that assignment to MIR reduces the likelihood of rearrest by 18.9 percentage points within the first six months, 18.4 percentage point within the first year, and 14.4 percentage points within the first two years. Relative to the control complier means, these effect sizes imply a 44 percent, 33 percent and 24 percent reduction in recidivism, respectively. The overall effect sizes hold up if we measure rearrest within three years of randomization (14.7 percentage points equivalent to 20 percent of the control-complier mean) and widens for four years (26.7 percentage points equivalent to 30 percent of the control-complier mean).

The final column models recidivism during the period 12 to 48 months post-randomization. This additional outcome allows us to assess whether referral to MIR has longer term effects on behavior beyond the immediate reduction in the likelihood of new arrest recidivism following the conference. Within 12 months, 99 percent of those assigned to MIR completed or failed the program. The average and median duration of MIR is 189 days. The estimate in Column (6) indicates that assignment to MIR reduces recidivism between years one and four post-randomization by 27 percentage points (equivalent to 37 percent of the control-complier mean). Regarding inference, the p-values from the one-tailed tests indicate significance at the five percent (or less) level of confidence for the six-month, 12-month, and 48-month outcomes, as well as for the outcome measuring recidivism one to four years after randomization. The p-values based on the randomization tests are generally larger, though we still observe significance at the five percent level for the six and 12-month outcomes and significance at the ten percent level for the 48-month and years one through four outcomes.

Turning to 2SLS estimates, the TOT effects are generally 1.3 to 1.4 times larger than the ITT effects. Relative to the control-complier means, the TOT estimates range from 30 percent of the recidivism rate for the 24-month outcome to 54 percent of recidivism occurring within six months and 51 percent of recidivism occurring in years one through four. The p-values of the TOT estimates are significant at the five percent level for all of the outcomes except the 24-month outcome (with a p-value of 0.0759) and the 36 month outcome (with a p-value of 0.0979). The p-values from the randomization tests align with the regular p-values.

Finally, we discuss whether non-compliance (being assigned to MIR and not enrolling) is related to recidivism propensities. Appendix Figure A.2 presents the cumulative failure functions for youth assigned to the control group and those assigned to the treatment group but did not enroll in MIR (the “never-takers”). The two curves are similar to each other, and a test for equality of the two curves fails to reject the null hypothesis that the two are equivalent ($p = 0.9485$). Moreover, the average rearrest rates at the bottom of Table 3 show that for all the outcomes, the control group

in which the youth is observed based on their date of randomization.

members and control compliers are quite close to one another, suggesting that non-compliance is unrelated to recidivism propensities.

5.1 External Validity of the Estimated Effects

As discussed above, youth who are eligible to participate in the RCT are a non-random sample (e.g., individuals affiliated with a gang are not eligible). However, one measure for external validity, arguably the most important one, is recidivism rates. Figure 3 presents various comparisons of the cumulative failure function for the MIR control group against the broader population of youth referred to JPD for felony offenses. Here we present two cumulative failure functions for the all youth sample: first, the cumulative failure function for all juvenile felony referrals and second, the cumulative failure function after re-weighting all juvenile felony referrals to match the youth in the MIR RCT in their predicted risk.²² Interestingly, the empirical failure function for all JPD referrals without re-weighting looks quite comparable to the empirical failure function for the control group. While the two groups have similar failure functions, the experimental sample have lower predicted likelihood of recidivism. Thus, as expected, re-weighting all JPD referrals to match the RCT in predicted recidivism yields a cumulative failure function that is somewhat below the empirical failure function for the control group at all points. These findings indicate that the MIR experimental sample might be negatively selected relative to the general population along certain dimensions. Regardless, we cannot reject the hypothesis that the failure function for the control group is equal to the unweighted cumulative failure function for all JPD felony referrals (for the weighted failure function the differences are significant at the 10 percent level). Our results suggest that the effects estimates presented above are likely valid for the broader population of youth who are typically charged with felonies in San Francisco.

5.2 The Effect of Enrollment vs. Completion of MIR

Our focus so far was on the ITT effect estimates of assignment to MIR and 2SLS estimates that re-scale the ITT by the share of individuals who enrolled into the program (Equations (2) and (3)). This section discusses an alternative 2SLS specification that re-scales the ITT based on the share of youths who completed the MIR program rather than the share who enrolled in it.

²²Specifically, we estimate a Probit model where the dependent variable is assignment to MIR and the regressors are predicted recidivism covariate indices for being arrested on a new offense within 6, 12, 18, 24, 30, 36, 42, and 48 months from when randomization took place (Section 4 and Appendix B describes the construction of the predicted recidivism scores based on individual characteristics and auxiliary observational data). Next we re-weight the cumulative failure function for all JPD felony referrals using the calculated propensity score and the re-weighting

procedure proposed by [Dinardo et al. \(1996\)](#). Specifically, $\text{weight}_i = \begin{cases} \frac{\Pr(E_i=1|X_i)}{1-\Pr(E_i=1|X_i)} \cdot \frac{1-\Pr(E_i=1)}{\Pr(E_i=1)} & E_i = 0 \\ 1 & E_i = 1 \end{cases}$, where

E_i denotes whether individual i is in the MIR experimental sample or in the observational data on all felony referrals.

According to the principles of restorative justice, meeting with the harmed individual and completing the contract agreed with them can be transformative for the youth. A key question is whether completion of the program is needed for it to have an effect. The summary statistics in Figure 1 suggest that the ITT recidivism reducing effect of assignment to MIR is driven by the individuals who completed the program. Specifically, youth assigned to MIR but who did not enroll have rearrest rates within one year of 57.9 percent. This is similar to the 57.7 percent rate of those who enrolled but did not complete the program. These rates are also similar to the 56.8 percent rearrest rate among the control units. However, among youth who completed MIR, the rearrest rate is much lower (19.2 percent). These summary statistics suggest that it is plausible that program completion is particularly important.

Table 4 reports 2SLS estimates of the TOT using an indicator for completing MIR as the treatment of interest. The TOT effects are larger than when using enrollment since the first stage (i.e., effect on completion) is smaller. Columns (1) and (2) report that completion of MIR reduces the likelihood of a rearrest within six and 12 months by about 36 percentage points. Relative to the control-complier means, these are reductions of 76 and 65 percent. The impacts on rearrests along longer time horizons are also large and range from 38 to 56 percent.

What factors lead individuals who enrolled to complete the MIR program? Appendix Table A.2 compares the characteristics of individuals who completed the program relative to those who did not (conditional on enrolling in MIR). A few patterns are noteworthy. First, in terms of the youths' demographic characteristics, sex and age do not predict completion. However, Black youth are less likely to complete the program than non-Black (mainly Hispanic) youth. Specifically, individuals who completed the program are 35 percentage points less likely to be Black. Second, there is suggestive evidence that individuals who are less likely to recidivate are completing the program with higher likelihoods. They have lower predicted recidivism, fewer prior arrests, and were older at their first arrest. Third, while youth who completed the program are generally less likely to recidivate based on observables, they can still be individuals charged with serious offenses. For example, youth arrested for an assault offense (the more severe type of crime in our sample) are more likely to complete the program than not.²³

5.3 Robustness Checks

5.3.1 Covariate Adjustment

To test the robustness of our difference-in-means comparisons to any finite sample imbalances in covariates, we also present ITT and TOT effects estimates that adjust for any imbalances in the predicted likelihood of a future arrest. To limit researcher degrees of freedom in deciding how to adjust for covariates (e.g., which controls to include in the model or not), we pre-specified

²³Individuals who participated in a restorative justice conference almost always completed the program. Only three individuals out of the 28 who did not complete the program participated in a restorative justice conference.

our procedure for conducting covariate adjustment with an eye on parsimony. The challenge is that including all the covariates in the regression model can meaningfully reduce the degrees of freedom since there can be many relevant variables.²⁴ Instead, we follow an alternative approach of adjusting only for a single summary index that can be viewed as a dimension reduction of the relevant information from all the covariates into a single factor. Specifically, we calculate the predicted likelihood of a future arrest using data on all the youth referred to JPD for felony offenses who did not participate in the MIR RCT (i.e., Column (4) of Table 2). We use this auxiliary data to estimate a simple OLS model where the dependent variable is an indicator of a future arrest for a new offense and the controls include sex, race/ethnicity, number of past arrests, number of past felony arrests, fixed effects for the category of the referral offense, age at arrest, and age at first arrest. We then use the estimated model coefficients (estimated only using the auxiliary data) and the covariates of study participants to generate the predicted likelihood of a future arrest and add this summary index to our base specification. Appendix Figure B.1 shows that the predicted future arrest index is highly correlated with the likelihood of a new arrest observed in the experimental sample, and we cannot reject a coefficient of one. In Appendix B, we present a more detailed description of this covariate adjustment procedure.²⁵

Appendix Table A.3 presents the results. The table is structured similarly to Table 3 but also includes the coefficient on the predicted recidivism index, which is a weighted average of the pre-treatment covariates. The weights are determined by the degree to which each covariate is predictive of the outcome in the auxiliary observational data on the full population of juveniles arrested for a felony offense in San Francisco. While there are small changes in coefficients, the results are very close to those from the models omitting the predicted recidivism index control variable. Thus, any finite imbalances in covariates do not impact our treatment effect estimates.

5.3.2 Differences in Effects Across Cohorts

The analyses above are based on individuals who have been assigned to MIR between late 2013 and May 2019. Our long-run (i.e., four-year) effects on recidivism are estimated using the earlier cohorts for which we have a longer time horizon to measure rearrests. While there is variation in the time horizon that we observe a youth for post-randomization, for all the youth in our sample we observe recidivism within at least one year from randomization. Next, we focus on rearrests occurring within six months or one year, for which we have a balanced sample, and examine

²⁴For example, controlling for the broad category of the current offense, age and other demographics, number of past offenses, number of past felony offenses, and the types of past offenses can result in a model with more than 30 covariates. Adding 30-plus covariates to the model in Equation (1) will have a meaningful impact on the number of degrees of freedom. To avoid needing to choose which covariates to include in the model, we opted for the option of controlling for a single index that summarizes the predictive ability of all the covariates as we describe in detail in Appendix B.

²⁵This is the procedure used to calculate the predicted likelihood of new arrest averages presented at the bottom of Table 2 and Appendix Table A.2. The idea of using auxiliary observational data to improve the accuracy of experimental estimates has been proposed in other studies (e.g., Gagnon-Bartsch et al., 2020).

differences across cohorts. Appendix Table A.4 reports 2SLS estimates with and without cohort fixed effects. The estimated effects are almost identical with and without the cohort fixed effects. Thus, in our balance sample the effects of MIR are similar across cohorts.

Lastly, in Appendix Table A.4, the coefficients on the cohort indicators are negative, indicating that, on average, the rearrest rate is lower in the more recent cohorts. In Appendix Figure A.3, we compare the rearrest rates of youth assigned to MIR in different time horizons (cohorts) and the control group. The figure shows that the later cohorts assigned to MIR generally have lower rearrest rates than the earlier cohorts. Thus, our findings suggest that the long-run effects would be at least as large if we observed four-year rearrest rates for all of our samples and not only for the earlier cohorts.

5.3.3 Robustness to Including or Not Including Arrests Due to Probation Violations

Our main measure of recidivism includes all rearrests including those that are the result of probation violations. As MIR can impact the case processing outcomes, using measures that include technical probation violation might classify similar behavior by the youth as recidivism depending on whether the individual is on probation. Moreover, whether one is on probation is impacted by MIR given the higher conviction rate for control group members. Therefore, to show our results are robust to the way recidivism is measured, we also report effect estimates using only arrests that are caused by a new criminal incident. Appendix Figure A.4 shows that when including only arrests for new criminal incidents, MIR has large and statistically significant recidivism-reducing effects that are similar to those documented in Figure 2. Moreover, the 2SLS and ITT estimates of the impacts of MIR (Appendix Table A.5) are similar to those reported in Table 3. These analyses confirm that the estimated effects of MIR are not sensitive to the decision of whether or not to measure recidivism using all rearrests or using only rearrests for new criminal incidents.

5.3.4 The COVID-19 Period

In this section we present a key robustness check pertaining to the overlap of the MIR observation period and the COVID epidemic. Recall that study subjects are randomized into MIR through October 2019 and our observations period extends through the end of 2020. Hence, many of the youth in the study have later observation periods that overlap with a time period when Californians (and Bay Area residents in particular) were subject to stringent stay-at-home orders.

Appendix Figure A.5 presents Kaplan-Meier estimates of the failure functions by treatment group where we have truncated the observation period for all youth to end at March 15, 2020.²⁶ Note, this truncation causes us to lose sample, especially for the observation periods beyond 16 months post-initial arrest. Nonetheless, the patterns we observe here are similar to what we observe

²⁶On March 16, 2020 San Francisco and five other Bay Area counties enacted a strict shelter-in-place orders that greatly reduced social interactions outside of the home and closed all in-person instruction in public schools throughout the region.

for the failure functions using un-truncated observation periods. Large disparities in the failure functions open up soon after the initial arrest and persist throughout the observation period. Both inference strategies reject the null hypothesis of equal failure functions for the treatment and control groups when we focus on arrests for any new offenses. For rearrests that result in a conviction, treatment group members are still less likely to be convicted for a new offense relative to the youth in the control group. Moreover, similar patterns also emerge when examining effects on more severe interactions such as a rearrest for a new felony offense or rearrests for new offenses that are as severe as the original offense.

6 Mechanisms

In this section, we discuss the potential mechanisms by which MIR may be impacting recidivism. One hypothesis is that restorative justice conferencing is a transformative event that causes youth to change their behavior. Furthermore, the meeting with the victim and the support the youth receives throughout the process from the conference coordinator and the agreement monitor can cause a lasting change in the youth’s behavior. For example, one of the program participants describes “In conference, we talked about how the person was harmed and how I harmed myself. I started to see things differently when I had to write a letter to myself about all the harm I had done.” Moreover, interviews and discussions with CWW staff (e.g., conference coordinator) reinforce the hypothesis that restorative justice conferencing is an emotionally impactful event for the youth that can potentially cause lasting effects on their behavior.

In addition to being a restorative justice conferencing program, MIR is also a diversion program. It diverts youth after the prosecutor decided to file charges but before charges have been filed. As a diversion program MIR influences case outcomes. Table 5 compares the conviction and case outcomes of individuals assigned to MIR relative to the control group. Panel (a) shows that among youth assigned to MIR, only five percent of the cases result in a felony conviction relative to 20 percent in the control group. Thus, assignment to MIR reduces the likelihood of being convicted of a felony offense by 15 percentage points. However, unlike criminal records in the adult system, in California, employers generally cannot use juvenile arrests and proceedings (e.g., charges, convictions) in employment decisions (Assembly Bill No. 1843) and this information does not appear in criminal history background checks. This restriction limits the degree to which interactions with the criminal justice system as juveniles can impact one’s labor market options.

Previous studies found that among adults a felony conviction can increase recidivism and negatively impact labor market outcomes (e.g., [Mueller-Smith and Schnepel, 2020](#); [Augustine et al., 2021](#)). Next, we examine whether the impacts of MIR can be explained by the case not resulting in a felony conviction. To do this, we restrict the sample only to individuals whose case did not result in a felony conviction and among them compare the rearrest rate of youth assigned to MIR relative to the control group. Importantly, restricting the sample only to cases that did not

result in a felony conviction biases us, if anything, against finding any effects of MIR. Individuals whose case results in a felony conviction are likely of higher risk of recidivism (at least not of lower risk), and 20 percent of the control group will have a felony conviction relative to only five percent of the treatment group.²⁷ Thus, removing these cases drops more individuals of potentially higher recidivism risk from the control group. Figure 4 shows that even among this sub-sample, youth participants assigned to MIR had a meaningfully lower likelihood of being rearrested both in the first few months as well as after four years from randomization. Moreover, the differences in rearrest rates between the two groups are statistically significant.

The results in Figure 4 show that even with this sample restriction to cases without a felony conviction, the MIR participants still have lower rearrest rates, indicating that the effects of MIR likely operate through its restorative justice conferencing component.²⁸ Furthermore, the ITT effects of MIR on the likelihood of being rearrested range from 14 to 27 percentage points, which is meaningfully larger than the effects of MIR on the case not resulting in a felony conviction. Thus, the primary channel for MIR’s recidivism-reducing effects is its conferencing component.

Another channel through which case outcomes can impact recidivism is the likelihood a youth will be removed from their home. However, Panel (b) of Table 5 shows that there is no difference in home removal rates among individuals assigned to MIR relative to the control group suggesting that differential home removal rates also cannot explain our findings. There is also no evidence of any type of differential “incapacitation” between the two groups.

Having established that restorative justice conferencing is the primary mediating channel for MIR’s effects, we next discuss the potential channels through which it may be altering behavior. One possibility is that the program changes an individual’s perspective, making them more aware of their impacts on others, increasing empathy for those on the receiving end of a criminal offense, and rendering the participant more deliberative, considerate, and self-reflective. Such an interpretation is consistent with the findings in past research that restorative justice conferencing appears to be most effective when the offense involves an identifiable victim (Sherman and Strang, 2007). Moreover, the shame associated with facing a crime victim during the conference might provide a more powerful deterrent than typical sanctions meted out through prosecution. Notably, a face-to-face meeting with the victim in which the youth listens to the harmful impacts of their actions on others is an emotional event that can influence an individual in multiple ways.

²⁷Appendix Figure A.6 shows that among all the youth arrested for a felony offense, those who are eventually convicted of a felony are meaningfully more likely to be rearrested in the future.

²⁸Agan et al. (2020) emphasize the importance of having any criminal history, even for non-violent misdemeanor offenses, on the likelihood of recidivism. MIR diverts youth before charging but after an arrest took place. Thus, both individuals in the control and treatment groups will have a criminal record of the initial arrest. In addition, the record of the arrest will still be visible to prosecutors in future cases (if the youth is rearrested). However, as we mentioned above, in California, juvenile arrests and proceedings (e.g., charges, convictions) do not appear in employment criminal history background checks.

7 Concluding Remarks

The MIR project randomized youth charged with select felonies to either prosecution as usual or referral to a restorative justice alternative. Take-up rates for those assigned to the MIR treatment were relatively high, especially compared to past restorative justice experiments. Assignment to MIR causes a large reduction in the likelihood of future arrest, both in the immediate aftermath of the initial arrest as well as up to four years following. We observe this for arrests overall and for more severe felony arrests.

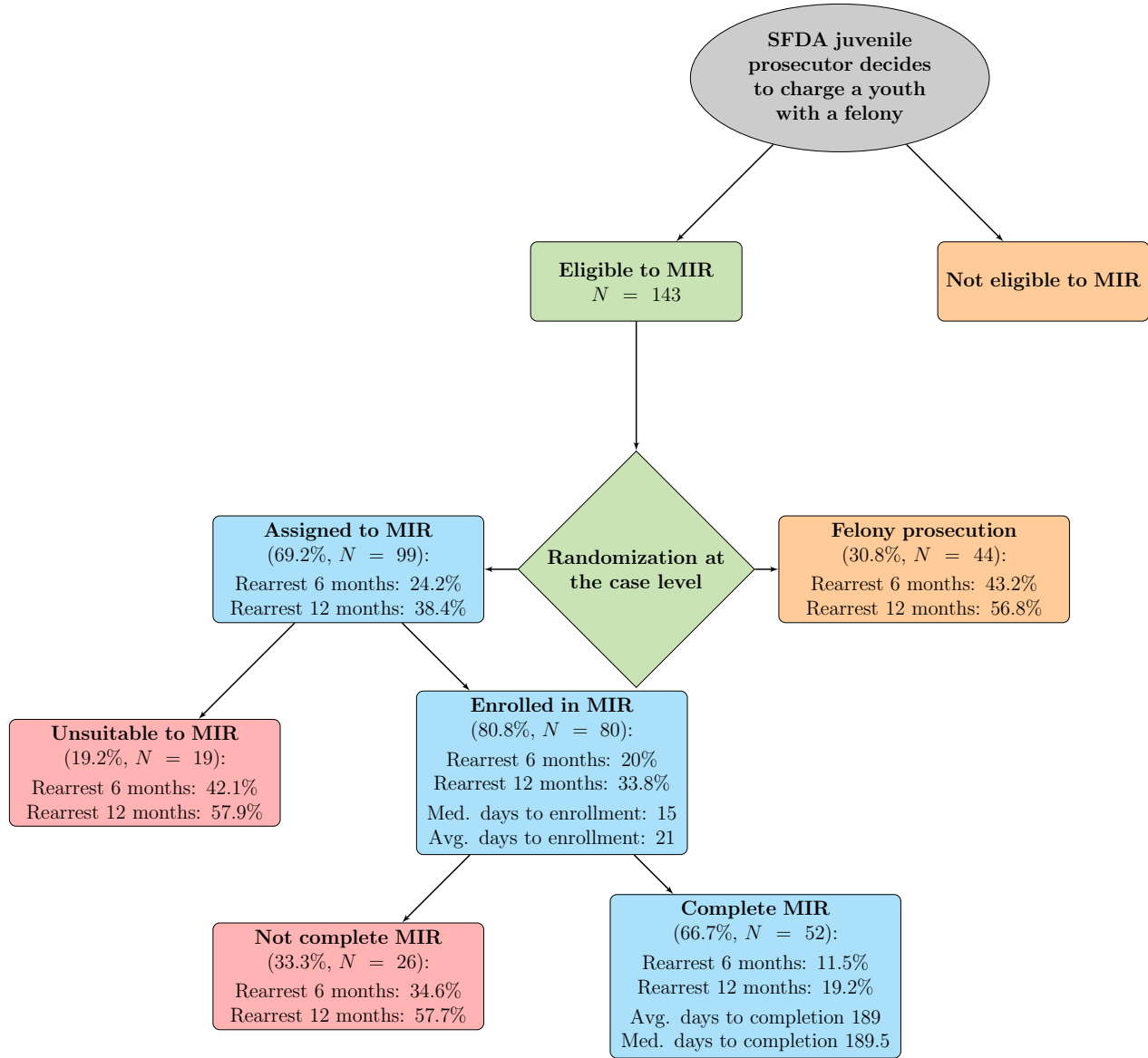
The large effects we observe for the MIR program raise questions concerning the mechanisms that are driving the results as well as aspects of the MIR program that may make this a particularly effective intervention. Beginning with the latter, there are two aspects of MIR that deserve further discussion. First, the MIR program treated youth charged with more serious offenses that are typically not eligible for similar restorative justice or diversion programs. It may be the case that there is simply more opportunity to reduce the likelihood of future arrests even for juveniles charged with serious offenses than generally understood or that interventions targeted at less serious offenses may simply widen the net of the criminal justice system and apply an invasive intervention to instances that do not merit such an intervention. Given the small scale of the experiment and the relatively narrow range of offenses qualifying for the program, we cannot evaluate whether effect size varies with risk of the individual (as defined by severity of controlling offense).

Second, unlike prior restorative justice experiments where youth assigned to the control group were funnelled into various alternative diversion programs, control group members under MIR faced felony prosecution. Moreover, treatment group members who declined to participate or who failed to successfully complete MIR faced felony prosecution. We suspect that the high take-up and compliance rates were driven by the deferred-adjudication structure of this intervention. However, we do not have experimental variation in this particular aspect to evaluate this hypothesis.

To conclude, our findings show that juvenile restorative justice conferencing can reduce recidivism among youth charged with relatively serious offenses and be an effective alternative to traditional criminal justice practices.

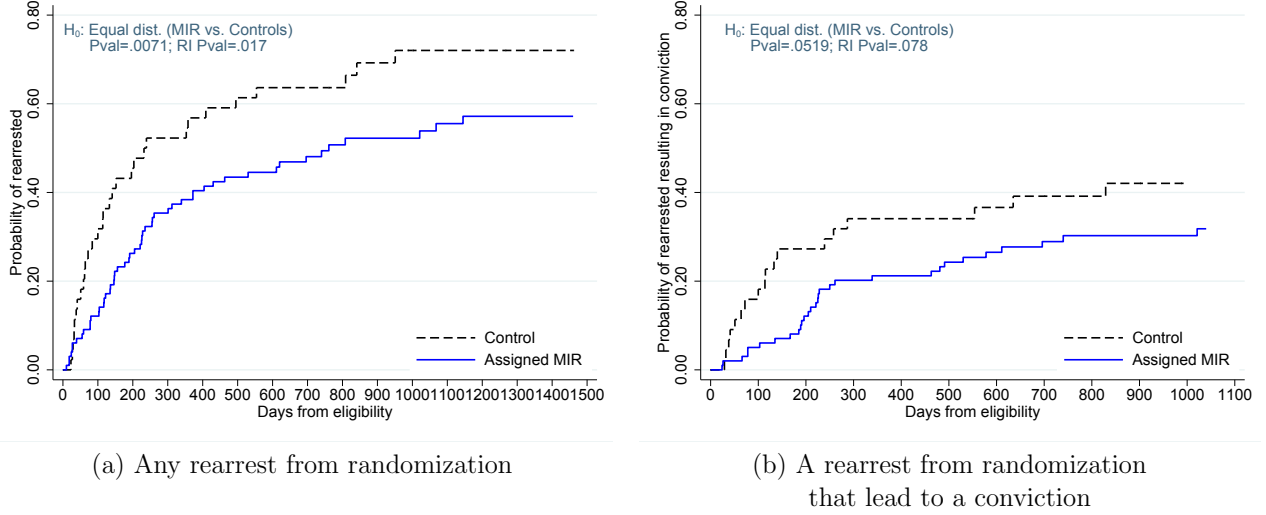
Figures

Figure 1: Make-it-Right Assignment, Enrollment, and Completion Process and Distribution



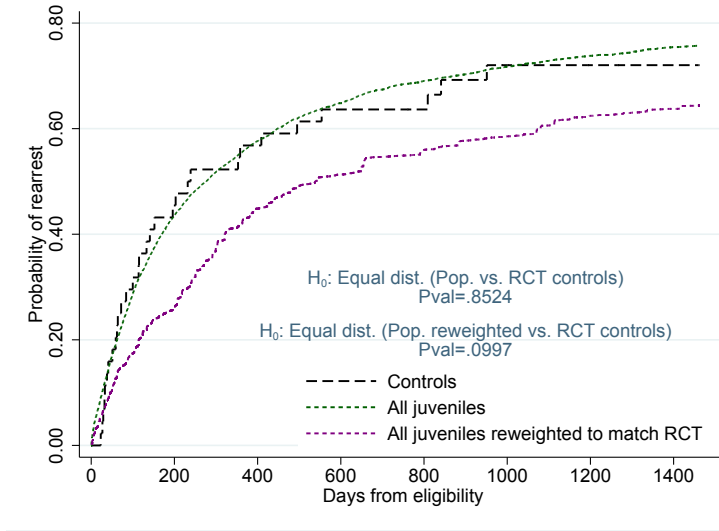
Notes: This figure depicts the process through which youths are assigned to the Make-it-Right (MIR) program. The process starts with the San Francisco District Attorney (SFDA) office receiving a reference / charges from the police or probation department regarding a youth. The SFDA juvenile prosecutor then decides whether or not to charge the youth with a felony offense. Our starting point (grey circle) is when the juvenile prosecutor decided to charge the youth with a felony offense. Then eligibility to MIR is assessed. If eligible, the youth is randomized to either MIR or to the control group which faces a traditional criminal prosecution for a felony offense. After being assigned to MIR the suitability of the youth to the program is being assessed by Community Works West (CWW) who administrates the MIR restorative justice conferencing program. Among the individuals assigned to MIR, 80.8% will be found suitable and continue to enroll in MIR and 19.2% will be determined as unsuitable and will face a standard criminal prosecution for a felony offense. Among the youths who enrolled into MIR, 66.7% will continue to complete the program and their charges will be permanently dropped upon completion. Another 33.3% will not succeed to complete the program and will face standard prosecution for their felony offense. The figure also reports the rearrest rates within six months and 12 months from the date in which youth are randomized into the different treatment regimes.

Figure 2: Recidivism Rates of Juveniles Randomly Assigned To Make-it-Right Relative to the Experimental Control Group



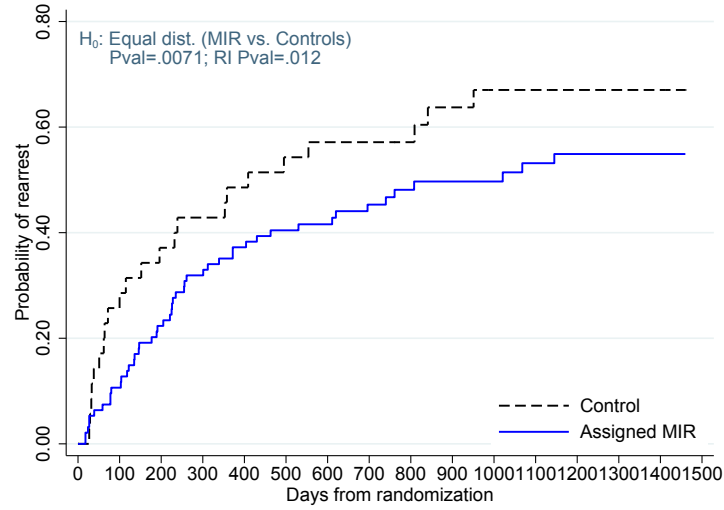
Notes: This figure plots Kaplan-Meier estimates of the failure function for being rearrested within four years from the date of randomization between assignment to Make-it-Right (MIR) and the control group which faces traditional criminal prosecution. Panel (a) plots Kaplan-Meier estimates for any rearrest and Panel (b) only counts rearrests that resulted in a conviction. In both plots we report p-values from an hypothesis test for whether the failure functions are the same among MIR participants and controls. We use the Peto–Peto–Prentice test for equality of failure functions (for a detailed description see [Klein and Moeschberger, 2006](#)). We report two types of p-values: “Pval” which is based on standard variance formulas and “RI Pval” which is based on randomization inference using 1,000 simulations in which we randomly assigned cases to MIR and controls and calculated the distribution of the test statistic under the null of no treatment effect. Both p-values are one-sided, rejecting the null hypothesis that MIR is not preferable to the control regime when the MIR Kaplan-Meier failure curve dominates that of the control units. We pre-specified in the pre-analysis plan that would to conduct only one-sided tests that reject the null when MIR is preferable to the control regime in reducing recidivism.

Figure 3: A Comparison Between the Make-it-Right Experimental Sample and the Full Population of Juveniles Arrested for a Felony Offense in San Francisco



Notes: This figure plots Kaplan-Meier estimates of the failure function of being rearrested. It compares the rearrest rates of the experimental control group (dashed black line) and the full populations of youth arrested for a felony offense in San Francisco (most of which are not eligible for MIR). The dotted green line reports Kaplan-Meier estimates of the rearrest rates of the full population of youth arrested for a felony offense in San Francisco and the dotted purple line reports re-weights estimates to make the full population similar to the experimental sample in terms of the predicted recidivism measures. The re-weighting is done using the procedure proposed by [Dinardo et al. \(1996\)](#). We report p-values from an hypothesis test for whether the failure functions are the same or not. We use the Peto–Peto–Prentice test for equality of failure functions (for a detailed description see [Klein and Moeschberger, 2006](#)). The p-values in these hypothesis tests are two-sided since we did not pre-specified in the pre-analysis plan how they will be conducted.

Figure 4: The Effects of Make-it-Right Among the Sub-Sample of Individuals Who Did Not Incur a Felony Conviction in Their Case Outcomes



Notes: This figure plots Kaplan-Meier estimates of the failure function for being rearrested within four years from the date of randomization into eligibility for Make-it-Right (MIR) or the control group which faces traditional criminal prosecution. In this figure (unlike Figure 2), we restrict the sample to individuals who case did not result in a felony conviction. As can be seen from the figure, even among this restricted population the individuals who have been assigned to MIR have meaningfully lower rearrest rates. We report p-values from an hypothesis test for whether the failure functions are the same among MIR participants and controls. We use the Peto–Peto–Prentice test for equality of failure functions (for a detailed description see [Klein and Moeschberger, 2006](#)). We report two types of p-values: “Pval” which is based on standard variance formulas and “RI Pval” which is based on randomization inference using 1,000 simulations in which we randomly assigned cases to MIR and controls and calculated the distribution of the test statistic under the null of no treatment effect. Both p-values are one-sided, rejecting the null hypothesis that MIR is not preferable to the control regime when the MIR Kaplan-Meier failure curve dominates that of the control units. We pre-specified in the pre-analysis plan that we would conduct only one-sided tests that reject the null when MIR is preferable to the control regime in reducing recidivism.

Tables

Table 1: Description of the Make-it-Right Restorative Justice Conferencing Program

Stage	Description/Examples of Activities
Suitability Assessment	CWW coordinator has initial meeting(s) with the youth (responsible party) and his/her family to determine if they are suitable for restorative justice community conferencing. The youth must agree to participate and demonstrate reflection and accountability to self, family, community, and person harmed. The main question that the coordinator asks when determining suitability is whether they feel confident putting the responsible youth in front of the person harmed. A youth who is unwilling to take responsibility will be deemed unsuitable at this point.
Pre-Conference (post-enrollment)	<p>The CWW coordinator holds several pre-conference meetings with the responsible youth and their family/support person to prepare them for the conference, typically between three to four. The youth must finalize an apology letter ahead of the conference. This is a reflective apology, for example, how do I feel about my actions now? If I had to do it over again, what would I do differently? What would I like the victim to know? What can or will I do to make up for what I did?</p> <p>The conference coordinator also conducts preparation meetings with the harmed party. These meeting aim to set expectations from the conference and for understanding the limitations the youth is facing.</p>
Conference	<p>The conference begins with the youth (responsible party) reading the apology letter to the harmed party.</p> <p>Next there is a roundtable discussion on how to address the four quadrants of the harm to: self, victim, family, and community.</p> <p>The conference results in a consensus-based plan of action (i.e., an agreement) for the youth's accountability and to prevent the youth from engaging in future criminal activity. The agreement's objective is to restore welfare by addressing the four quadrants of the harm: self, victim, family, and community.</p> <p>All parties sign the agreement. This concludes the formal involvement of the harmed party. Multiple conferences can be held until the plan is developed—if no plan is developed, the youth is referred back to the SFDA for prosecution.</p>
Examples of Agreement Activities	<p>Academic: Tutoring sessions, meet high school attendance requirements, make a plan for college/technical school application. Employment: Make a goal to create a resume and apply to a certain number of jobs. Reflection Writing: Journaling, poems, and/or essays reflecting on opportunities for self-improvement. Yoga: Attend a set number of yoga classes. Anger management: Attend/complete a set number of anger management sessions. Restitution: Identify amount to provide to harmed party and/or community. Chores: Keeping one's room clean, taking out the trash, helping with dinner, etc. Goal is to help the youth engage more in family life. Family Systems Therapy: Counseling sessions with identified family members. Community Service: Youth repairs harm done to the community by performing a set number of community service hours at a local organization.</p>
Agreement Implementation	<p>After the conference, the Huckleberry Youth agreement monitor debriefs with the youth. They finalize the details of the restorative plan (i.e., the agreement) and set target completion dates.</p> <p>Youth and agreement monitor meet on a weekly basis to review the youth's progress toward completion of the plan. The meetings are not only about making sure the youth is on track to finish the plan; they also discuss other issues that the youth is facing and develop a plan to address them.</p>

Notes: This table summarizes the activities and sequence of events of the Make-it-Right (MIR) program. The initials CWW refer to Community Works West which is a non-profit community organization that specializes in restorative justice conferencing and administrates the MIR program in San Francisco. The initials SFDA refer to the San Francisco District Attorney.

Table 2: Summary Statistics of Make-it-Right Experimental Sample and the Full Sample of Juveniles Arrested for a Felony Offense in San Francisco

	(1) MIR controls	(2) Assigned MIR	(3) MIR Compliers	(4) All juveniles
Demographics:				
Male	0.909	0.889 [0.756]	0.900	0.796
Black	0.500	0.531 [0.788]	0.487	0.608
Hispanic	0.318	0.323 [0.966]	0.359	0.238
Age	16.023	16.091 [0.814]	16.113	15.851
Criminal history:				
Any past arrests	0.273	0.273 [1.000]	0.250	0.565
Number of past arrests	0.705	0.374 [0.229]	0.313	1.903
Any past felony arrests	0.068	0.141 [0.184]	0.125	0.468
Number of past felony arrests	0.068	0.162 [0.134]	0.125	0.960
Age at first criminal offense	14.750	15.198 [0.091]	15.269	14.227
Type of most severe offense:				
Homicide/Manslaughter	0.000	0.000 [.]	0.000	0.019
Sex offense	0.000	0.000 [.]	0.000	0.013
Robbery	0.000	0.030 [0.087]	0.000	0.344
Assault	0.159	0.131 [0.690]	0.138	0.264
Burglary	0.318	0.434 [0.245]	0.487	0.140
Theft	0.636	0.657 [0.833]	0.713	0.236
Drug	0.000	0.000 [.]	0.000	0.071
Weapons	0.000	0.020 [0.165]	0.025	0.098
Other	0.205	0.293 [0.447]	0.287	0.472
Predicted recidivism:				
Pred. recidivism 6 months	0.305	0.296 [0.569]	0.289	0.421
Pred. recidivism 12 months	0.433	0.426 [0.654]	0.417	0.564
Pred. felony recidivism 6 months	0.182	0.185 [0.856]	0.179	0.254
Pred. felony recidivism 12 months	0.313	0.303 [0.504]	0.296	0.369
Joint F-test of MIR assignment on covariates p-value	0.757			
Joint F-test of MIR compliers and never-takers covariates p-value	0.840			
Number of observations	44	99	80	6272
Number of individuals	44	99	80	3333

Notes: The table reports summary statistics (means) of the individuals randomly assigned to the control group (Column (1)), to Make-it-Right (MIR) (Column (2)), the compliers—those assigned to MIR and who also enrolled into the program (Column (3)), and the full population of juveniles arrested for a felony offense between October 2010 and November 2020. The square parenthesis in Column (2) report p-values for whether the differences in each characteristic between Columns (1) and (2) are different than zero. The mean characteristics of compliers in Column (3) are calculated using the standard formula from [Abadie \(2003\)](#). Specifically, using a 2SLS regression of a covariate interacted with an indicator for enrollment into MIR (i.e., $MIR_i \cdot X_i$) as the outcome, an indicator for MIR enrollment as the endogenous treatment, and instrumenting using an indicator for whether the youth was randomly assigned to control or MIR. Note that not all individuals assigned to MIR took-up the program. The take-up rate is about 80 percent. The joint F-tests at the bottom of the table are based on randomization inference using 1,000 random placebo permutations of treatment assignment.

Table 3: The Effects of Assignment (ITT) to and Participation (TOT) in Make-it-Right on the Likelihood of Being Arrested in the Subsequent Four Years

	(1)	(2)	(3)	(4)	(5)	(6)
	6 months	12 months	24 months	36 months	48 months	12-48 months
Panel (a)	<i>2SLS</i>					
Participated in MIR (treated)	-0.234 (0.103) {0.0120} [0.0030]	-0.228 (0.111) {0.0211} [0.0050]	-0.184 (0.128) {0.0759} [0.0722]	-0.196 (0.151) {0.0979} [0.1150]	-0.363 (0.165) {0.0157} [0.0230]	-0.368 (0.199) {0.0344} [0.0250]
Panel (b)	<i>Reduced form</i>					
Assigned to MIR (ITT)	-0.189 (0.084) {0.0132} [0.0140]	-0.184 (0.092) {0.0237} [0.0410]	-0.144 (0.103) {0.0830} [0.1130]	-0.147 (0.118) {0.1092} [0.1680]	-0.267 (0.133) {0.0249} [0.0850]	-0.270 (0.154) {0.0423} [0.1040]
Panel (c)						
First-Stage coefficient	0.808 (.0463)	0.808 (.0463)	0.781 (.0558)	0.750 (.0676)	0.736 (.0832)	0.736 (.0832)
Rearrest rate among controls	0.432	0.568	0.632	0.750	0.833	0.667
Rearrest rate among compliers controls	0.434	0.566	0.606	0.745	0.876	0.726
Includes controls	No	No	No	No	No	No
Number of observations	143	143	120	100	71	71

Notes: The table reports estimates of the effects of Make-it-Right (MIR) on the likelihood of a future arrest. Each cell in the table reports four numbers: the point estimate, standard error clustered at the case level, a one-sided p-value using the cluster-robust standard errors, and a one-sided p-value using randomization inference based on 1,000 random permutations. The compliers rearrest rates under the control regime (bottom of the table) are calculated using the standard formulas from [Imbens and Rubin \(1997\)](#) and [Abadie \(2002\)](#). Specifically, using a 2SLS regression of the outcome interacted with an indicator for enrollment into MIR (i.e., $(1 - \text{MIR}_i) \cdot \text{Rearrest}_i$) as the outcome, an indicator for not enrolling into MIR (i.e., $(1 - \text{MIR}_i)$) as the endogenous treatment, and instrumenting using an indicator for whether the youth was randomly assigned to control or MIR. Note that not all individuals assigned to MIR took-up the program. The take-up rate is about 75% and is reported at the bottom of the table (i.e., the First-Stage coefficient). The number of observations changes across the columns because the sample in each of the regressions is restricted to individuals that are observed at least the mentioned time horizon (e.g., 48 months in Column (5)) after the date of randomization.

Table 4: The Effects of Completing Make-it-Right on the Likelihood of Being Arrested in the Subsequent Four Years

	(1)	(2)	(3)	(4)	(5)	(6)
	6 months	12 months	24 months	36 months	48 months	12-48 months
Panel (a)	<i>2SLS</i>					
Completed MIR (treated)	-0.361 (0.159) {0.0125} [0.0010]	-0.351 (0.166) {0.0178} [0.0101]	-0.268 (0.176) {0.0653} [0.1180]	-0.286 (0.205) {0.0837} [0.1690]	-0.525 (0.214) {0.0084} [0.0460]	-0.531 (0.276) {0.0295} [0.0380]
Panel (b)						
First-Stage coefficient	0.525 (.0726)	0.525 (.0726)	0.537 (.0825)	0.515 (.0974)	0.509 (.1231)	0.509 (.1231)
Rearrest rate among controls	0.432	0.568	0.632	0.750	0.833	0.667
Rearrest rate among compliers controls	0.476	0.543	0.541	0.743	0.932	0.827
Includes controls	No	No	No	No	No	No
Number of observations	143	143	120	100	71	71

Notes: The table reports estimates of the effects of Make-it-Right (MIR) completion on the likelihood of a future arrest. The key difference between the estimates in this table relative to Table 3 is that here the treatment is defined as *completion* rather than enrollment in MIR. Here the key assumption is that assignment to MIR impacts recidivism only by impacting the likelihood that an individual will complete the MIR program. In other words, participation in the program without completion has no effect on recidivism. Each cell in the table reports four numbers: the point estimate, standard error clustered at the case level, a one-sided p-value using the cluster-robust standard errors, and a one-sided p-value using randomization inference based on 1,000 random permutations. The compliers rearrest rates under the control regime (bottom of the table) are calculated using the standard formulas from Imbens and Rubin (1997) and Abadie (2002). Specifically, using a 2SLS regression of the outcome interacted with an indicator for enrollment into MIR (i.e., $(1 - \text{MIR}_i) \cdot \text{Rearrest}_i$) as the outcome, an indicator for not enrolling into MIR (i.e., $(1 - \text{MIR}_i)$) as the endogenous treatment, and instrumenting using an indicator for whether the youth was randomly assigned to control or MIR. Note that not all individuals assigned to MIR took-up the program. The take-up/completion rate is about 50 percent and is reported at the bottom of the table (i.e., the First-Stage coefficient). The number of observations changes across the columns because the sample in each of the regressions is restricted to individuals that are observed at least the mentioned time horizon (e.g., 48 months in Column (5)) after the date of randomization.

Table 5: A Comparison of Charging, Convictions, and Case Dispositions Between Individuals Assigned to Make-it-Right Relative to Those Assigned to the Control Group

	(1) Control	(2) Assigned MIR
Convicted offenses:		
Convicted of misdemeanor or felony	0.386	0.152
Convicted of misdemeanor	0.205	0.121
Convicted of felony	0.205	0.0505
Disposition (least to most severe):		
Charges dismissed	0.364	0.737
Informal probation	0.364	0.121
Probation	0.0455	0.0101
Wardship probation	0.136	0.0707
Out of home placement	0.0227	0.0202
Number of individuals	44	99

Notes: The table reports information on charging (i.e., whether an individual was charged with an offense or the prosecutor decided to drop the case), convictions, and dispositions for the individuals in our experimental sample. These case outcomes are for the criminal incidents that lead the individual to be in our sample. We do not include charging, conviction, or disposition information from future criminal incidents.

References

- Abadie, Alberto**, “Bootstrap Tests for Distributional Treatment Effects in Instrumental Variable Models,” *Journal of the American Statistical Association*, 2002, pp. 284–293.
- , “Semiparametric instrumental variable estimation of treatment response models,” *Journal of econometrics*, 2003, *113* (2), 231–263.
- Abrams, David S**, “Estimating the deterrent effect of incarceration using sentencing enhancements,” *American Economic Journal: Applied Economics*, 2012, *4* (4), 32–56.
- Agan, Amanda and Sonja Starr**, “Ban the Box, Criminal Records, and Racial Discrimination: A Field Experiment,” *The Quarterly Journal of Economics*, 2018, *133* (1), 191–235.
- , **Jennifer Doleac, and Anna Harvey**, *Misdemeanor Prosecution*, College Station, TX: Texas AM Working Paper, 2020.
- Aizer, Anna and Joseph J. Doyle**, “Juvenile Incarceration, Human Capital, and Future Crime: Evidence from Randomly Assigned Judges,” *The Quarterly Journal of Economics*, 2015, *130* (2), 759–803.
- Angel, Caroline M., Lawrence W. Sherman, Heather Strang, Barak Ariel, Sarah Bennett, Nova Inkpem, Anne Keane, and Therese S. Richmond**, “Short-Term Effects of Restorative Justice Conferences on Post-Traumatic Stress Symptoms Among Robbery and Burglary Victims: A Randomized Controlled Trial,” *Journal of Experimental Criminology*, 2014, *10*, 291–307.
- Angrist, Joshua D, Guido W Imbens, and Donald B Rubin**, “Identification of causal effects using instrumental variables,” *Journal of the American statistical Association*, 1996, *91* (434), 444–455.
- Arteaga, Carolina**, “Parental Incarceration and Children’s Educational Attainment,” 2020. Working paper.
- Augustine, Elsa, Johanna Lacoe, Alissa Skog, and Steven Raphael**, *The Impact of Felony Diversion in San Francisco*, Berkeley, CA: University of California Berkeley Working Paper, 2021.
- Becker, Gary S.**, “Crime and Punishment: An Economic Approach,” *Journal of Political Economy*, 1968, *76* (2), 169–217.
- Berg, Mark T., Eric A. Stewart, Christopher J. Schreck, and Ronald L. Simmons**, “The victim-offender overlap in context: Examining the role of neighborhood street culture,” *Criminology*, 2012, *50*, 359–390.
- Bhuller, Manudeep, Gordon B. Dahl, Katrine V. Løken, and Magne Mogstad**, “Intergenerational Effects of Incarceration,” *AEA Papers and Proceedings*, 2018, *108*, 234–40.
- , —, —, and —, “Incarceration, Recidivism, and Employment,” *Journal of Political Economy*, 2020, *128* (4), 1269–1324.

- Blattman, Christopher, Julian C Jamison, and Margaret Sheridan**, “Reducing crime and violence: Experimental evidence from cognitive behavioral therapy in Liberia,” *American Economic Review*, 2017, *107* (4), 1165–1206.
- Bloom, Howard S**, “Accounting for no-shows in experimental evaluation designs,” *Evaluation review*, 1984, *8* (2), 225–246.
- Bonta, James, Suzanne Wallace-Capretta, Jennifer Rooney, and Kevin Mcanoy**, “An Outcome Evaluation of Restorative Justice Alternative to Incarceration,” *Contemporary Justice Review*, 2002, *5* (4), 319–338.
- Braithwaite, John**, *Crime, Shame, and Reintegration*, Cambridge, U.K.: Cambridge University Press, 1989.
- Britto, Diogo, Paolo Pinotti, and Breno Sampaio**, “The effect of job loss and unemployment insurance on crime in Brazil,” 2020.
- Brooks, Alison**, *Moving forward: Two approaches to repairing the harm through restorative justice*, American University, 2013. PhD dissertation.
- Bushway, Shawn D.**, “Labor Market Effects of Permitting Employer Access to Criminal History Records,” *Journal of Contemporary Criminal Justice*, 2004, *20* (3), 276–291.
- Christensen, Garret and Edward Miguel**, “Transparency, reproducibility, and the credibility of economics research,” *Journal of Economic Literature*, 2018, *56* (3), 920–80.
- Chung, EunYi and Joseph P Romano**, “Exact and asymptotically robust permutation tests,” *The Annals of Statistics*, 2013, *41* (2), 484–507.
- Cuellar, Allison E., Larkin S. McReynolds, and Gail A. Wasserman**, “A Cure For Crime: Can Mental Health Treatment Diversion Reduce Crime Among Youth,” *Journal of Policy Analysis and Management*, 2006, *25* (1), 197–214.
- Damm, Anna Piil and Christian Dustmann**, “Does growing up in a high crime neighborhood affect youth criminal behavior?,” *American Economic Review*, 2014, *104* (6), 1806–32.
- Davidson, Janet, George King, Jens Ludwig, and Steven Raphael**, “Managing Pretrial Misconduct: An Experimental Evaluation of HOPE Pretrial.”, Technical Report, GSPP Working Paper 2019.
- Dinardo, J, NM Fortin, and T Lemieux**, “Labor market institutions and the distribution of Wages, 1973-1992: A semiparametric approach,” *Econometrica*, 1996, *64* (5), 1001–1044.
- Dobbie, Will, Hans Grönqvist, Susan Niknami, Mårten Palme, and Mikael Priks**, “The intergenerational effects of parental incarceration,” Technical Report, National Bureau of Economic Research 2018.
- Doleac, Jennifer L and Benjamin Hansen**, “The unintended consequences of “ban the box”: Statistical discrimination and employment outcomes when criminal histories are hidden,” *Journal of Labor Economics*, 2020, *38* (2), 321–374.

- Draca, Mirko, Theodore Koutmeridis, and Stephen Machin**, “The changing returns to crime: do criminals respond to prices?,” *The Review of Economic Studies*, 2019, 86 (3), 1228–1257.
- Drago, Francesco, Roberto Galbiati, and Pietro Vertova**, “The Deterrent Effects of Prison: Evidence from a Natural Experiment,” *Journal of Political Economy*, 2009, 117 (2), 257–280.
- Eren, Ozkan and Naci Mocan**, “Juvenile Punishment, High School Graduation, and Adult Crime: Evidence from Idiosyncratic Judge Harshness,” *Review of Economics and Statistics*, 2021, 103 (1), 34–47.
- Fisher, RA et al.**, “The design of experiments,” *The design of experiments.*, 1935, (1nd Ed).
- Fulkerson, Andrew**, “The use of victim impact panels in domestic violence cases: A Restorative Justice Approach,” *Contemporary Justice Review*, 2001, 4, 355–368.
- Gagnon-Bartsch, Johann A, Adam C Sales, Edward Wu, Anthony F Botelho, Luke W Miratrix, and Neil T Heffernan**, “Precise Unbiased Estimation in Randomized Experiments using Auxiliary Observational Data,” 2020.
- Heller, Sara B, Anuj K Shah, Jonathan Guryan, Jens Ludwig, Sendhil Mullainathan, and Harold A Pollack**, “Thinking, fast and slow? Some field experiments to reduce crime and dropout in Chicago,” *The Quarterly Journal of Economics*, 10 2017, 132 (1), 1–54.
- Hirschi, Travis and Michael Gottfredson**, “Age and the explanation of crime,” *American journal of sociology*, 1983, 89 (3), 552–584.
- Hjalmarsson, Randi**, “Juvenile Jails: A Path to the Straight and Narrow or to Hardened Criminality?,” *Journal of Law and Economics*, 2009, 52 (4), 779–809.
- Huntington, Margaret, Lisa Quan, Sarah Riley, and Richard Zarrella**, *Make It Right: An Evaluation of a Youth Restorative Justice Program in San Francisco*, University of California, Berkeley: Goldman School of Public Policy, 2017.
- Huttunen, K., M. Kaila, M. Kaila, T. Kosonen, and E. Nix**, “Shared Punishment? The Impact of Incarcerating Fathers on Child Outcomes,” 2019. Working paper.
- i Vidal, Jordi Blanes and Giovanni Mastrobuoni**, “Police patrols and crime,” 2018.
- Imbens, Guido and Donald Rubin**, “Estimating Outcome Distributions for Compliers in Instrumental Variables Models,” *The Review of Economic Studies*, 1997, 64 (4), 555–574.
- Imbens, Guido W. and Joshua D. Angrist**, “Identification and Estimation of Local Average Treatment Effects,” *Econometrica*, 1994, 62 (2), 467–475.
- Jeong, Seokjin, Edmund F McGarrell, and Natalie Kroovand Hipple**, “Long-term impact of family group conferences on re-offending: The Indianapolis restorative justice experiment,” *Journal of Experimental Criminology*, 2012, 8 (4), 369–385.
- Katz, Lawrence F, Jeffrey R Kling, and Jeffrey B Liebman**, “Moving to opportunity in Boston: Early results of a randomized mobility experiment,” *The Quarterly Journal of Economics*, 2001, 116 (2), 607–654.

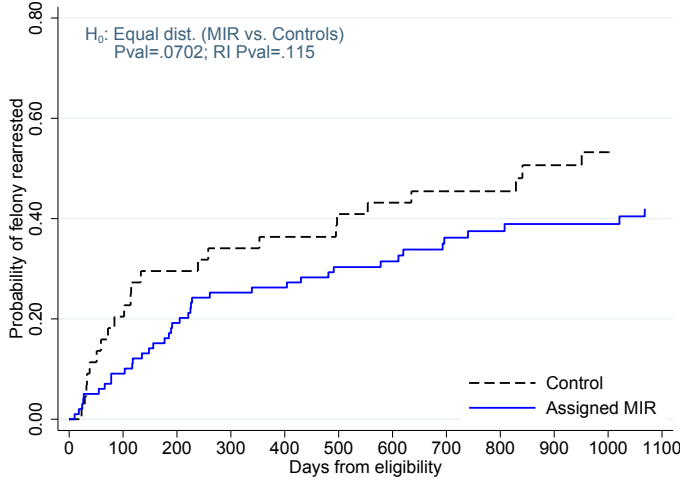
- Khanna, Gaurav, Carlos Medina, Anant Nyshadham, Christian Posso, and Jorge Tamayo**, “Job Loss, Credit, and Crime in Colombia,” *American Economic Review: Insights*, 2021, 3 (1), 97–114.
- Kirchmaier, Tom, Stephen Machin, Matteo Sandi, and Robert Witt**, “Prices, policing and policy: the dynamics of crime booms and busts,” *Journal of the European Economic Association*, 2020, 18 (2), 1040–1077.
- Klein, John P and Melvin L Moeschberger**, *Survival analysis: techniques for censored and truncated data*, Springer Science & Business Media, 2006.
- Kling, Jeffrey R, Jens Ludwig, and Lawrence F Katz**, “Neighborhood effects on crime for female and male youth: Evidence from a randomized housing voucher experiment,” *The Quarterly Journal of Economics*, 2005, 120 (1), 87–130.
- Kuziemko, Ilyana**, “How should inmates be released from prison? An assessment of parole versus fixed-sentence regimes,” *The Quarterly Journal of Economics*, 2013, 128 (1), 371–424.
- Li, Xinran and Peng Ding**, “General forms of finite population central limit theorems with applications to causal inference,” *Journal of the American Statistical Association*, 2017, 112 (520), 1759–1769.
- Little, Simon, Anne Stewart, and Nicole Ryan**, “Restorative Justice Conferencing: Not a Panacea for the Overrepresentation of Australia’s Indigenous Youth in the Criminal Justice System,” *International Journal of Offender Therapy and Comparative Criminology*, 2018, 62 (13), 4067–4090.
- Livingstone, Nuala, Geraldine Macdonald, and Nicola Carr**, “Restorative justice conferencing for reducing recidivism in young offenders (aged 7 to 21),” *Cochrane Library*, 2013, (2).
- Lofstrom, Magnus and Steven Raphael**, “Crime, the Criminal Justice System, and Socioeconomic Inequality,” *The Journal of Economic Perspectives*, 2016, 30 (2), 103–126.
- Ludwig, Jens, Greg J Duncan, and Paul Hirschfield**, “Urban poverty and juvenile crime: Evidence from a randomized housing-mobility experiment,” *The Quarterly Journal of Economics*, 2001, 116 (2), 655–679.
- McCold, Paul and Benjamin Wachtel**, *Restorative Policing Experiment: The Bethlehem Pennsylvania Police Family Group Conferencing Project*, Pipersville, PA: Community Service Foundation, 1998.
- McGarrell, Edmund F**, *Restorative Justice Conferences as an Early Response to Young Offenders*, Washington, D.C.: Office of Juvenile Justice and Delinquency Prevention, U.S. Department of Justice, 2001.
- **and Natalie Kroovand Hipple**, “Family group conferencing and re-offending among first-time juvenile offenders: The Indianapolis experiment,” *Justice Quarterly*, 2007, 24 (2), 221–246.
- Ministry of Justice**, *Restorative Justice: Best Practice in New Zealand*, Wellington, New Zealand: Ministry of Justice, 2004.

- Mitchell, Ojmarrh, David B. Wilson, Amy Eggers, and Doris L. MacKenzie**, “Assessing the Effectiveness of Drug Courts on Recidivism: A Meta-Analytical Review of Traditional and Non-Traditional Drug Courts,” *Journal of Criminal Justice*, 2012, 40 (1), 60–71.
- Mueller-Smith, Michael and Kevin T. Schnepel**, “Diversion in the Criminal Justice System,” *Review of Economic Studies*, 2020, rdaa030, <https://doi.org/10.1093/restud/rdaa030>.
- National Research Council**, *The Growth of Incarceration in the United States: Exploring Causes and Consequences*, Washington, D.C.: National Academies Press, 2014.
- Neal, Derek and Armin Rick**, “The prison boom and the lack of black progress after Smith and Welch,” Technical Report, National Bureau of Economic Research 2014.
- Norris, Samuel, Matthew Pecenco, and Jeffrey Weave**, “The Effects of Parental and Sibling Incarceration: Evidence from Ohio,” 2020. Working paper.
- Olken, Benjamin A**, “Promises and perils of pre-analysis plans,” *Journal of Economic Perspectives*, 2015, 29 (3), 61–80.
- Owens, Emily, Aria Golestani, and Kerri Raissian**, “Specialization in Criminal Courts: Decision Making, Recidivism, and Re-victimization in Domestic Violence Courts in Tennessee,” 2021. Working paper.
- Raphael, Steven**, “The New Scarlet Letter? Negotiating the U.S. Labor Market with a Criminal Record,” 2014. Report.
- and **Michael A. Stoll**, *Why Are So Many Americans in Prison*, New York, NY: Russell Sage Foundation, 2013.
- Rose, E**, “Who gets a second chance? effectiveness and equity in supervision of criminal offenders,” *The Quarterly Journal of Economics*, 2020.
- Rose, Evan and Yotam Shem-Tov**, “How does incarceration affect reoffending? Estimating the dose-response function,” *Journal of Political Economy*, Forthcoming.
- Rose, Evan K**, “Does banning the box help ex-offenders get jobs? Evaluating the effects of a prominent example,” *Journal of Labor Economics*, 2021, 39 (1), 79–113.
- Sampson, Robert J. and Janet L. Lauritsen**, “Deviant lifestyles, proximity to crime and the offender- victim link in personal violence,” *Journal of Research on Crime and Delinquency*, 1990, 27, 110–139.
- Schnepel, Kevin T**, “Good jobs and recidivism,” *The Economic Journal*, 2018, 128 (608), 447–469.
- Shapland, Joanna, Anne Atkinson, Helen Atkinson, James Dignan, Lucy Edwards, Jeremy Hibbert, Marie Howes, Jennifer Johnstone, Gwen Robinson, and Angela Sorsby**, “Does Restorative Justice Affect Reconviction: The Fourth Report from the Evaluation of Three Scheme,” Technical Report, Ministry of Justice Research Series 10/08, United Kingdom 2008.

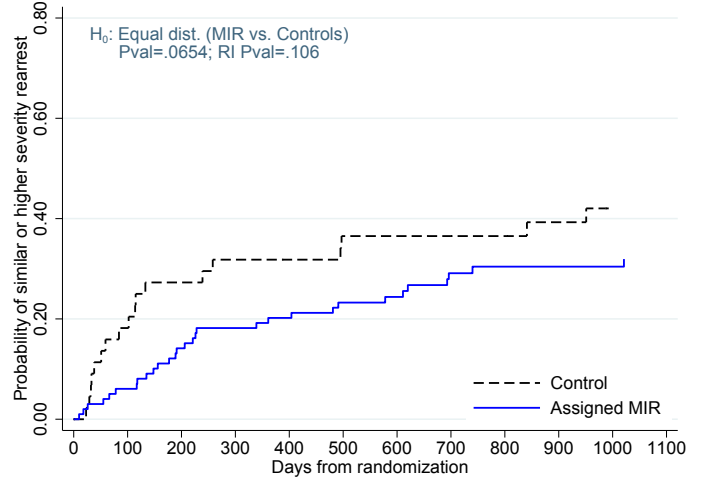
- Sherman, Lawrence and Heather Strang**, “Restorative Justice: The Evidence,” Technical Report, The Smith Institute, London, United Kingdom 2007.
- Sherman, Lawrence W., Heather Strang, Geoffrey Barnes, Daniel J. Woods, Sarah Bennett, Nova Inkpen, Dorothy Newbury-Birch, Meredith Rossner, Caroline Angel, Malcolm Mearns, and Molly Slothower**, “Twelve Experiments in Restorative Justice: the Jerry Lee Program of Randomized Trials of Restorative Justice Conferences,” *Journal of Experimental Criminology*, 2015, *11*, 501–540.
- Strang, Heather, John Braithwaite et al.**, *Restorative justice and family violence*, Cambridge University Press, 2002.
- , **Lawrence W Sherman, Evan Mayo-Wilson, Daniel Woods, and Barak Ariel**, “Restorative justice conferencing (RJC) using face-to-face meetings of offenders and victims: Effects on offender recidivism and victim satisfaction. A systematic review,” *Campbell Systematic Reviews*, 2013, *9* (1), 1–59.
- Tella, Rafael Di and Ernesto Schargrodsky**, “Do police reduce crime? Estimates using the allocation of police forces after a terrorist attack,” *American Economic Review*, 2004, *94* (1), 115–133.
- Ulmer, Jeffrey Todd and Darrell J Steffensmeier**, “The age and crime relationship: Social variation, social explanations,” in “The nurture versus biosocial debate in criminology: On the origins of criminal behavior and criminality,” SAGE Publications Inc., 2014, pp. 377–396.
- Umbreit, Mark S. and Marilyn Peterson Armour**, “Restorative Justice and Dialogues: Impact, Opportunities, and Challenges in the Global Community,” *Washington University Journal of Law Policy*, 2011, *36*, 65–89.
- Wilson, David B, Ajima Olaghere, and Catherine S Kimbrell**, *Effectiveness of restorative justice principles in juvenile justice: A meta-analysis* 2018.
- Wu, Jason and Peng Ding**, “Randomization tests for weak null hypotheses in randomized experiments,” *Journal of the American Statistical Association*, 2020, pp. 1–16.
- Yang, Crystal S**, “Local labor markets and criminal recidivism,” *Journal of Public Economics*, 2017, *147*, 16–29.
- Young, Alwyn**, “Channeling fisher: Randomization tests and the statistical insignificance of seemingly significant experimental results,” *The Quarterly Journal of Economics*, 2019, *134* (2), 557–598.
- Zhao, Anqi and Peng Ding**, “Covariate-adjusted Fisher randomization tests for the average treatment effect,” *Journal of Econometrics*, 2021.

A Additional Figures and Tables

Figure A.1: A Comparison of Recidivism Rates Between Juveniles Randomly Assigned To Make-it-Right Relative to the Experimental Control Group Along Different Measures of the Severity of Reoffending



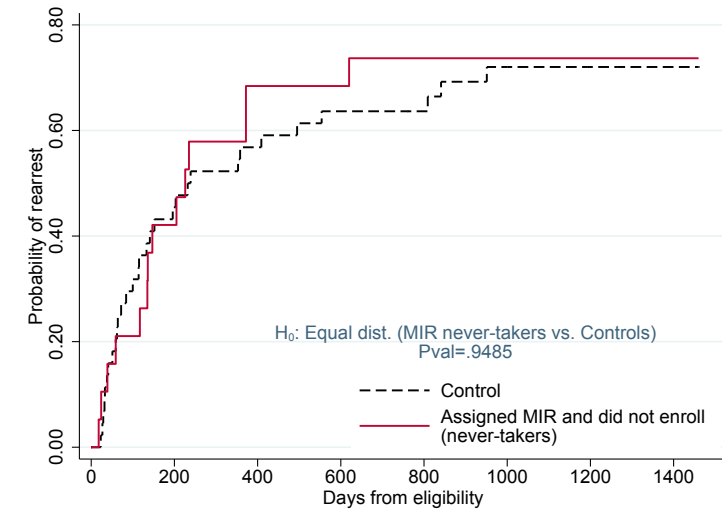
(a) Any felony rearrest from randomization



(b) Any rearrest from randomization for offenses at least as severe as the original offense

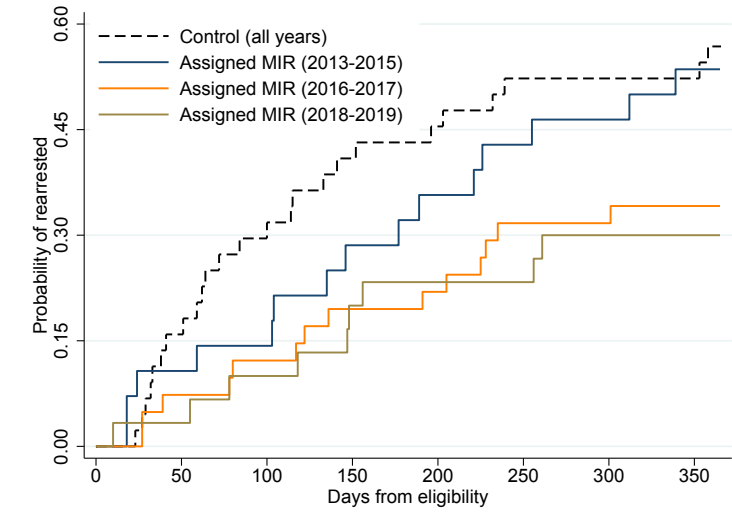
Notes: This figure plots Kaplan-Meier estimates of the failure function for being rearrested within four years from the date of randomization between assignment to Make-it-Right (MIR) and the control group which faces traditional criminal prosecution. Panel (a) plots Kaplan-Meier estimates for any rearrest and Panel (b) only counts rearrests that resulted in a conviction. In both plots we report p-values from an hypothesis test for whether the failure functions are the same among MIR participants and controls. We use the Peto–Peto–Prentice test for equality of failure functions (for a detailed description see [Klein and Moeschberger, 2006](#)). We report two types of p-values: “Pval” which is based on standard variance formulas and “PI Pval” which is based on permutation inference using 1,000 simulations in which we randomly assigned cases to MIR and controls and calculated the distribution of the test statistic under the null of no treatment effect. Both p-values are one-sided, rejecting the null hypothesis that MIR is not preferable to the control regime when the MIR Kaplan-Meier failure curve dominates that of the control units. We pre-specified in the pre-analysis plan that we would conduct only one-sided tests that reject the null when MIR is preferable to the control regime in reducing recidivism.

Figure A.2: A Comparison of Recidivism Rates Between the Control Group and Individuals Assigned to Make-it-Right But Who Did Not Enroll (Never-Takers)



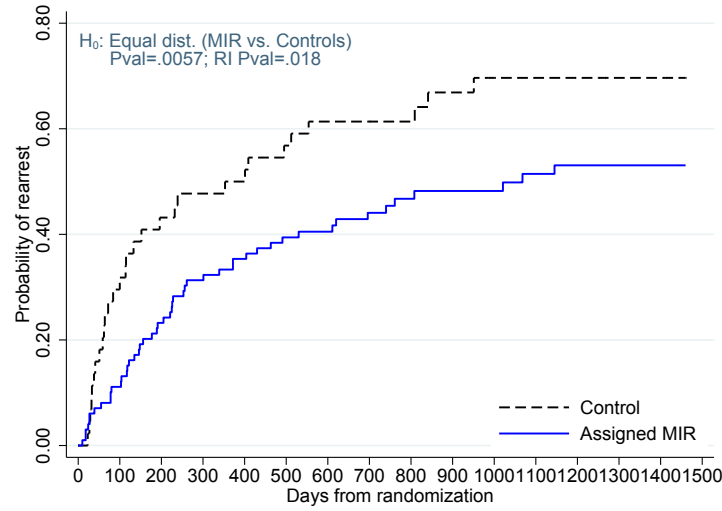
Notes: This figure plots Kaplan-Meier estimates of the failure function of being rearrested. The figure compares between the failure curves of the experimental control group (dashed black line) and of youth randomly assigned to Make-it-Right (MIR) but who did not enroll into the MIR program, i.e., the “never-takers” (solid red line). The p-value is from an hypothesis test for whether the failure functions are the same or not. We use the Peto–Peto–Prentice test for equality of failure functions (for a detailed description see [Klein and Moeschberger, 2006](#)). The p-values in these hypothesis test are two-sided since we did not pre-specify in the pre-analysis plan how they would be conducted.

Figure A.3: Rearrest Rates of Juveniles Assigned To Make-it-Right in Different Time Periods Relative to the Experimental Control Group



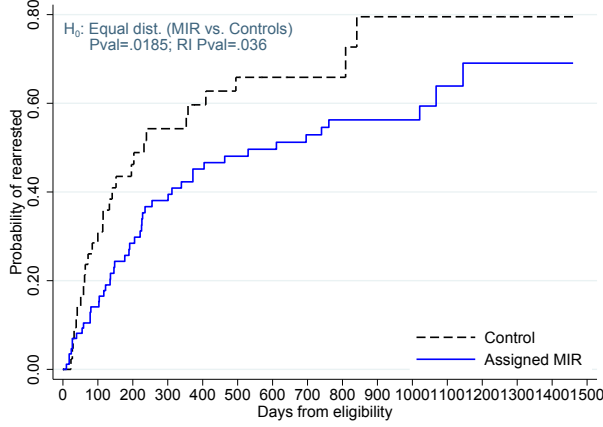
Notes: This figure plots Kaplan-Meier estimates of the failure function for being rearrested within one year from the date of randomization into eligibility for Make-it-Right (MIR) or the control group which faces traditional criminal prosecution. We observe at least one year post-randomization for all the individuals in our sample. Thus, when comparing the rearrest rates within one year we do not need to drop any observations. However, when examining the likelihood of a rearrest in a longer time horizon the sample decreases.

Figure A.4: Recidivism Rates of Juveniles Randomly Assigned To Make-it-Right Relative to the Experimental Control Group Using Only Rearrests for New Criminal Incidents

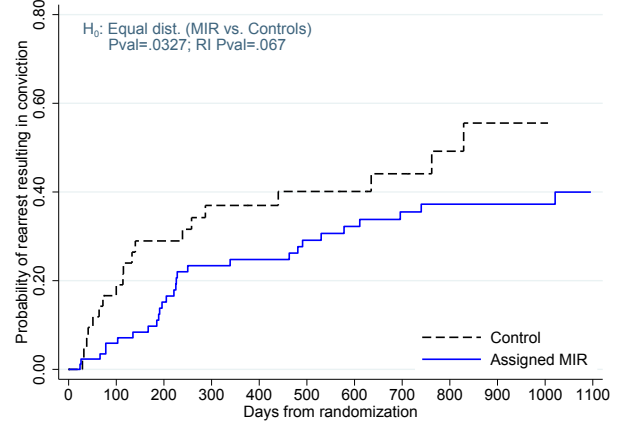


Notes: This figure plots Kaplan-Meier estimates of the failure function for being rearrested within four years from the date of randomization between assignment to Make-it-Right (MIR) and the control group which faces traditional criminal prosecution. Recidivism is measured using only rearrests for new criminal incidents. Specifically, rearrests for probation or warrants violations will not be included in this recidivism measure. We report p-values from an hypothesis test for whether the failure functions are the same among MIR participants and controls. We use the Peto–Peto–Prentice test for equality of failure functions (for a detailed description see [Klein and Moeschberger, 2006](#)). We report two types of p-values: “Pval” which is based on standard variance formulas and “RI Pval” which is based on randomization inference using 1,000 simulations in which we randomly assigned cases to MIR and controls and calculated the distribution of the test statistic under the null of no treatment effect. Both p-values are one-sided, rejecting the null hypothesis that MIR is not preferable to the control regime when the MIR Kaplan-Meier failure curve dominates that of the control units. We pre-specified in the pre-analysis plan that we would conduct only one-sided tests that reject the null when MIR is preferable to the control regime in reducing recidivism.

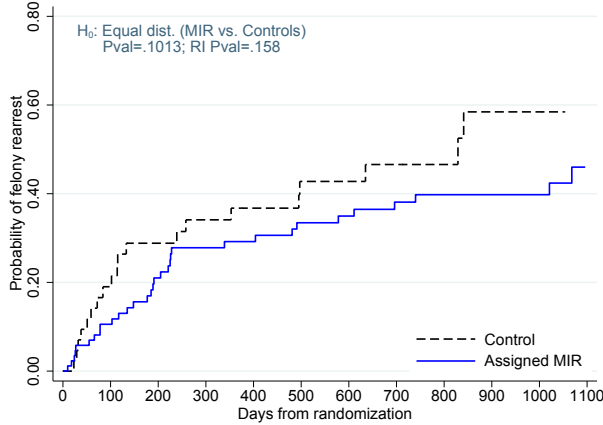
Figure A.5: Rearrest Rates of Juveniles Randomly Assigned to Make-it-Right Relative to the Experimental Control Group When Not Including Any Reoffending That Took Place After March 15, 2020 Which Is the Beginning of COVID-19 Restrictions in California



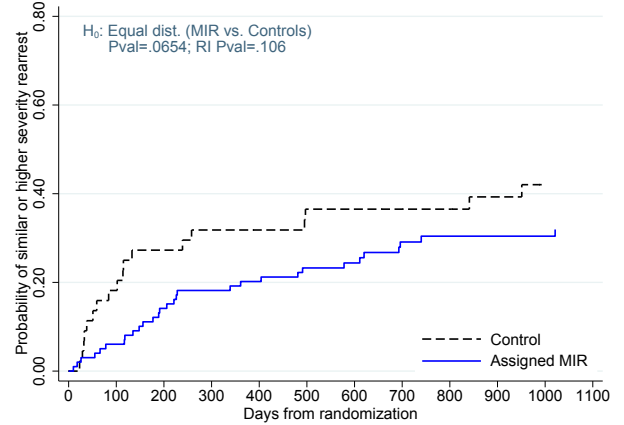
(a) Any rearrest from randomization



(b) A rearrest from randomization that leads to a conviction



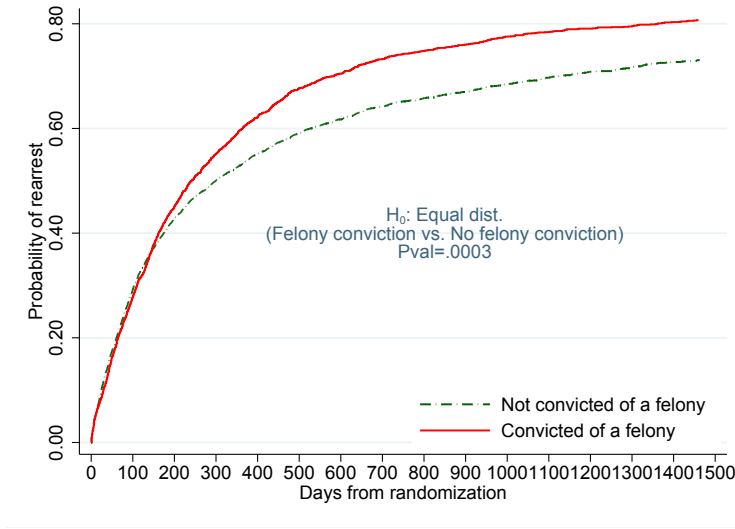
(c) Any felony rearrest from randomization



(d) Any rearrest from randomization for offenses at least as severe as the original offense

Notes: This figure plots Kaplan-Meier estimates of the failure function for being rearrested within four years from the date of randomization into eligibility for Make-it-Right (MIR) or the control group which faces regular criminal prosecution. Panel (a) plots Kaplan-Meier estimates for any rearrest and Panel (b) only counts rearrests that lead to a conviction. Panel (c) plots rearrests for felony offenses. Lastly, Panel (d) presents failure function including only rearrests for offenses that are as severe as the original offense. In all plots we report p-values from an hypothesis test for whether the failure functions are the same among MIR participants and controls. We use the Peto–Peto–Prentice test for equality of failure functions (for a detailed description see [Klein and Moeschberger, 2006](#)). We report two types of p-values: “Pval” which is based on standard variance formulas and “RI Pval” which is based on randomization inference using 1,000 simulations in which we randomly assigned cases to placebo MIR and controls and calculated the distribution of the test statistic under the null of no treatment effect. Both p-values are one-sided, rejecting the null hypothesis that MIR is not preferable to the control regime when the MIR Kaplan-Meier estimates of the failure curve dominates that of the control units. We pre-specified in the pre-analysis plan that we would conduct only one-sided tests that reject the null when MIR is preferable to the control regime in reducing recidivism.

Figure A.6: A Comparison of Recidivism Rates in the Full Population of Juveniles Arrested for a Felony Offense in San Francisco Between Individuals Whose Case Resulted in a Felony Conviction Relative to Those Who Did Not



Notes: This figure plots Kaplan-Meier estimates of the failure function for being rearrested within four years from the date the offense took place. The figure plots two failure function curves, one for individuals whose case resulted in a felony conviction (solid red line) and another for the individuals whose case did not result in a felony conviction. We report p-values from an hypothesis test for whether the failure functions are the same or not. We use the Peto-Peto-Prentice test for equality of failure functions (for a detailed description see [Klein and Moeschberger, 2006](#)). The p-values in these hypothesis tests are two-sided since we did not pre-specified in the pre-analysis plan how they will be conducted.

Table A.1: Summary Statistics of the Demographic Composition of the Victim (Harmed Party) and the Youth (Responsible Party) in the Make-it-Right Experimental Sample

	(1) Victim (harmed party)	(2) Youth (responsible party)
Age	35.60	16.09
Sex:		
Male	0.585	0.889
Victim and youth of same sex	0.523	.
Missing sex	0.343	0
Race/ethnicity:		
Black	0.0820	0.531
Hispanic	0.148	0.323
White	0.443	0.0729
Asian	0.328	0.0938
Victim and youth of same race	0.213	.
Missing race	0.384	0

Notes: The table reports summary statistics (means) of the demographic characteristics of the youth (responsible party) who have been assigned to Make-it-Right (MIR) and the victim (harmed party) of the related criminal incidents. The demographic information for the youth comes from the administrative records provided to use by the San Francisco District Office and the San Francisco Department of Juvenile Probation. The demographic information for the victims comes from Community Works West, which is the non-profit organization that implements MIR. As a result, we observe demographic information for only a subset of the victims.

Table A.2: Summary Statistics Comparing Individuals Who Completed Make-it-Right to Those Who Did Not Complete Make-it-Right Among the Youth Who Enrolled Into the Program

	(1)	(2)
	Did not complete MIR	Completed MIR
Demographics:		
Male	0.893	0.904
Black	0.714	0.360
Hispanic	0.286	0.400
Age	15.86	16.25
Criminal history:		
Any past arrests	0.321	0.212
Age at first criminal offense	14.93	15.46
Assault	0.107	0.154
Burglary	0.429	0.519
Theft	0.714	0.712
Weapons	0.0357	0.0192
Other	0.107	0.385
Predicted recidivism:		
Pred. recidivism 6 months	0.323	0.271
Pred. recidivism 12 months	0.450	0.400
Pred. felony recidivism 6 months	0.208	0.163
Pred. felony recidivism 12 months	0.329	0.278
Number of individuals	28	52

Notes: The table reports summary statistics (means) of the characteristics of individuals who enrolled into the to the Make-it-Right (MIR) program and did not complete it (Column (1)) and of individuals who enrolled and completed the program (Column (2)).

Table A.3: The Effects of Assignment (ITT) to and Participation (TOT) in MIR on the Likelihood of Being Arrested in the Subsequent Four Years When Including Controls

	(1)	(2)	(3)	(4)	(5)	(6)
	6 months	12 months	24 months	36 months	48 months	12-48 months
Panel (a)	<i>2SLS</i>					
Participated in MIR (treated)	-0.231 (0.104) {0.0137} [0.0000]	-0.223 (0.111) {0.0235} [0.0110]	-0.177 (0.125) {0.0792} [0.0480]	-0.191 (0.144) {0.0946} [0.0900]	-0.358 (0.165) {0.0167} [0.0200]	-0.363 (0.198) {0.0353} [0.0180]
Predicted Y(0)	0.270 (0.457)	0.528 (0.435)	0.652 (0.369)	0.798 (0.400)	0.560 (0.424)	0.392 (0.391)
Panel (b)	<i>Reduced form</i>					
Assigned to MIR (ITT)	-0.185 (0.085) {0.0151} [0.0090]	-0.179 (0.091) {0.0262} [0.0380]	-0.137 (0.100) {0.0860} [0.1240]	-0.143 (0.113) {0.1050} [0.1590]	-0.262 (0.132) {0.0259} [0.0710]	-0.267 (0.153) {0.0436} [0.0830]
Predicted Y(0)	0.408 (0.485)	0.661 (0.451)	0.746 (0.369)	0.912 (0.402)	0.743 (0.433)	0.496 (0.395)
Panel (c)						
First-Stage coefficient	0.802 (.0456)	0.803 (.0454)	0.776 (.0551)	0.747 (.0658)	0.732 (.0815)	0.734 (.0824)
Rearrest rate among controls	0.432	0.568	0.632	0.750	0.833	0.667
Rearrest rate among compliers controls	0.434	0.566	0.606	0.745	0.876	0.726
Includes controls	Yes	Yes	Yes	Yes	Yes	Yes
Number of observations	143	143	120	100	71	71

Notes: The table reports estimates of the effects of Make-It-Right (MIR) on the likelihood of a future arrest. Predicted $Y_i(0)$ is a summary index of the covariates based on how predictive they are on the outcome in the auxiliary observational data, in which none of the individuals were assigned to MIR. The construction of predicted $Y(0)$ based on the covariates is described in Section 4 and more formally in Appendix B. The only difference between this table and Table 3 is the inclusion of predicted $Y_i(0)$ in the regression specifications. Each cell in the table reports four numbers: the point estimate, standard error clustered at the case level, a one-sided p-value using the cluster-robust standard errors, and a one-sided p-value using randomization inference based on 1,000 random permutations. The compliers rearrest rates under the control regime (bottom of the table) are calculated using the standard formulas from Imbens and Rubin (1997) and Abadie (2002). Specifically, using a 2SLS regression of the outcome interacted with an indicator for enrollment into MIR (i.e., $(1 - \text{MIR}_i) \cdot \text{Rearrest}_i$) as the outcome, an indicator for not enrolling into MIR (i.e., $(1 - \text{MIR}_i)$) as the endogenous treatment, and instrumenting using an indicator for whether the youth was randomly assigned to control or MIR. Note that not all individuals assigned to MIR took-up the program. The take-up rate is about 75% and is reported at the bottom of the table (i.e., the First-Stage coefficient). The number of observations changes across the columns because the sample in each of the regressions is restricted to individuals that are observed at least the mentioned time horizon (e.g., 48 months in Column (5)) after the date of randomization.

Table A.4: 2SLS Estimates of the Effects of Enrollment Into Make-it-Right on Rearrest Within Six Months and One Year for a Balanced Sample With and Without Cohort Fixed Effects

	(1)	(2)	(3)	(4)
	6 months	6 months	12 months	12 months
Participated in MIR (treated)	-0.234 (0.103) {0.0120} [0.0000]	-0.233 (0.104) {0.0134} [0.0020]	-0.228 (0.111) {0.0211} [0.0080]	-0.212 (0.111) {0.0292} [0.0120]
2016-2017 cohort		-0.059 (0.088)		-0.159 (0.096)
2018-2019 cohort		-0.024 (0.099)		-0.212 (0.104)
First-Stage coefficient	0.808 (.0463)	0.798 (.0509)	0.808 (.0463)	0.798 (.0509)
Rearrest rate among controls	0.432	0.432	0.568	0.568
Rearrest rate among compliers controls	0.434	0.434	0.566	0.566
Includes controls				
Number of observations	143	143	143	143

Notes: The table reports estimates of the effects of enrollment into Make-It-Right (MIR) on the likelihood of a future arrest. Importantly, the sample size does not change across columns as we restricted attention to shorter time horizons in which a balanced sample is observed. Each cell in the table reports four numbers: the point estimate, standard error clustered at the case level, a one-sided p-value using the cluster-robust standard errors, and a one-sided p-value using randomization inference based on 1,000 random permutations. In Columns (1) and (3), no covariates are included in the model. In Columns (2) and (4), fixed effects for the time period/cohort in which the incident took place. The omitted category is the 2013-2015 cohort, i.e., the cases for which we observe the longest follow-up period. The compliers rearrest rates under the control regime (bottom of the table) are calculated using the standard formulas from [Imbens and Rubin \(1997\)](#) and [Abadie \(2002\)](#). Specifically, using a 2SLS regression of the outcome interacted with an indicator for enrollment into MIR (i.e., $(1 - \text{MIR}_i) \cdot \text{Rearrest}_i$) as the outcome, an indicator for not enrolling into MIR (i.e., $(1 - \text{MIR}_i)$) as the endogenous treatment, and instrumenting using an indicator for whether the youth was randomly assigned to control or MIR. Note that not all individuals assigned to MIR took-up the program. The take-up rate is about 75% and is reported at the bottom of the table (i.e., the First-Stage coefficient). The take-up rate is not affected by whether or not cohort fixed effects are included in the model.

Table A.5: The Effects of Assignment (ITT) to and Participation (TOT) in Make-it-Right on the Likelihood of Being Arrested for a New Criminal Incident in the Subsequent Four Years

	(1)	(2)	(3)	(4)	(5)	(6)
	6 months	12 months	24 months	36 months	48 months	12-48 months
Panel (a)	<i>2SLS</i>					
Participated in MIR (treated)	-0.244 (0.102) {0.0094} [0.0000]	-0.206 (0.113) {0.0351} [0.0040]	-0.197 (0.127) {0.0620} [0.0401]	-0.213 (0.150) {0.0791} [0.0840]	-0.440 (0.158) {0.0035} [0.0020]	-0.368 (0.199) {0.0344} [0.0250]
Panel (b)	<i>Reduced form</i>					
Assigned to MIR (ITT)	-0.197 (0.083) {0.0097} [0.0100]	-0.167 (0.093) {0.0373} [0.0520]	-0.154 (0.102) {0.0681} [0.0910]	-0.160 (0.118) {0.0892} [0.1360]	-0.324 (0.127) {0.0069} [0.0310]	-0.270 (0.154) {0.0423} [0.1040]
Panel (c)						
First-Stage coefficient	0.808 (.0463)	0.808 (.0463)	0.781 (.0558)	0.750 (.0676)	0.736 (.0832)	0.736 (.0832)
Rearrest rate among controls	0.409	0.500	0.605	0.719	0.833	0.667
Rearrest rate among compliers controls	0.444	0.519	0.588	0.723	0.902	0.726
Includes controls	No	No	No	No	No	No
Number of observations	143	143	120	100	71	71

Notes: The table reports estimates of the effects of Make-it-Right (MIR) on the likelihood of a future arrest. Each cell in the table reports four number: the point estimate, standard error clustered at the case level, a one-sided p-value using the cluster-robust standard errors, and a one-sided p-value using randomization inference based on 1,000 random permutations. The only difference between this table and Table 3 is that here only rearrests for new criminal incidents are used. Specifically, rearrests for probation or warrant violations will not be included. The compliers rearrest rates under the control regime (bottom of the table) are calculated using the standard formulas from Imbens and Rubin (1997) and Abadie (2002). Specifically, using a 2SLS regression of the outcome interacted with an indicator for enrollment into MIR (i.e., $(1 - \text{MIR}_i) \cdot \text{Rearrest}_i$) as the outcome, an indicator for not enrolling into MIR (i.e., $(1 - \text{MIR}_i)$) as the endogenous treatment, and instrumenting using an indicator for whether the youth was randomly assigned to control or MIR. Note that not all individuals assigned to MIR took-up the program. The take-up rate is about 75% and is reported at the bottom of the table (i.e., the First-Stage coefficient). The number of observations changes across the columns because the sample in each of the regressions is restricted to individuals that are observed at least the mentioned time horizon (e.g., 48 months in Column (5)) after the date of randomization.

B Covariates Adjustment in an RCT Using Auxiliary Observational Data

We observe two samples. The first is an experimental sample in which individuals were randomly assigned to either Make-it-Right (MIR) or the control group. Denote by Z_i assignment to MIR (the treatment of interest), let Y_i^s be the outcome of interest, and denote by X_i^s a vector of pre-treatment covariates. The second sample is much larger, however, the treatment (Z_i) was not assigned to any of the individuals in this sample. Denote by Y_i^p and X_i^p the outcome and observable characteristics of individuals in the larger auxiliary sample. In our setting, this sample includes all juveniles not eligible for MIR who have been arrested for a felony offense between October 2010 and November 2020.

In the experimental sample, assignment to MIR is done at random, however, due to the small number of observations there can still be some imbalances between the observable and unobservable characteristics of the individuals who were assigned to the treatment and control groups. Specifically, the bias term can be expressed as:

$$\begin{aligned} \mathbb{E}[Y_i^s|Z_i = 1] - \mathbb{E}[Y_i^s|Z_i = 0] = \\ \mathbb{E}[Y_i^s(1) - Y_i^s(0)|Z_i = 1] + \mathbb{E}[Y_i^z(0)|Z_i = 1] - \mathbb{E}[Y_i^s(0)|Z_i = 0] \end{aligned} \tag{B.1}$$

It is clear from Equation (B.1) that if we observed $Y_i^s(0)$ among both the control and treated units then we could control for it and correct for any potential finite sample imbalances.

It is common practice to use Ordinary-Least-Square (OLS) regression and estimate the following specification:

$$Y_i^s = \alpha Z_i + X_i^{s'} \beta^s + e_i^s \tag{B.2}$$

this model corrects for potential imbalances in observables between the treatment and control groups and with flexible/saturated enough controls it is completely non-parametric (i.e., it does not require making any parametric assumption such as linearity of the conditional expectation

function). Another motivation for controlling for X_i^s is increasing the statistical precision by improving the explanatory power of the OLS model.

The big challenge in Equation (B.2) is that increasing the dimensionality of X , i.e., including a greater number of covariates, entails a reduction in the number of degrees of freedom. Reducing the degrees of freedom can be costly in experiments with small sample sizes as is the case in our setting of the MIR program. Moreover, including only a subset of the potential X s raises the question of which covariates to include and adds another “researcher degree of freedom”. To overcome this problem, we use the auxiliary data on all juveniles arrested for a felony offense between October 2010 and November 2020 in San Francisco.

We begin by using auxiliary observational data to derive an estimator of $\widehat{X_i^{s'}\beta^s}$. In both the experimental and observational samples, we observe the same vector of observable characteristics. We estimate the following OLS specification in the auxiliary data:

$$Y_i^p = X_i^{p'}\beta^p + e_i^p \quad (\text{B.3})$$

next we use the estimated coefficient $\hat{\beta}^p$ to form our estimator of $X_i^{s'\beta^s}$:

$$\hat{Y}_{0i} \equiv X_i^{s'}\hat{\beta}^p \quad (\text{B.4})$$

and now we can estimate the following model:

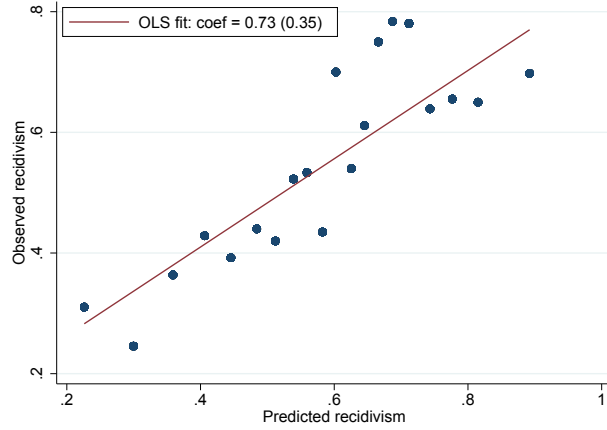
$$Y_i^s = \alpha Z_i + \gamma \hat{Y}_{0i} + \nu_i^s \quad (\text{B.5})$$

where $\nu_i^s = (\gamma \hat{Y}_{0i} - X_i^{s'}\beta^s) + e_i^s$. Note that, if $\hat{\beta}^p \approx \beta^s$, then $\nu_i^s = e_i^s$ and the specification in Equation (B.5) yields the same results as the one in Equation (B.2) while using only one degree of freedom since only a *single* covariate is included in the model.

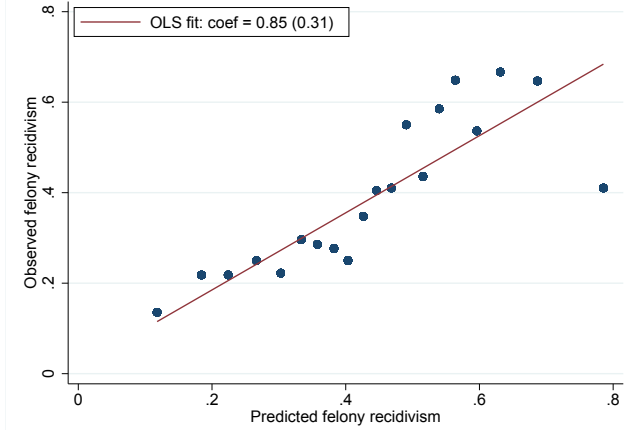
To validate that our predicted recidivism index ($X_i^{s'}\hat{\beta}^p$) is predictive of recidivism in the experimental sample, we examine the relationship between Y_i^s and $X_i^{s'}\hat{\beta}^p$. Figure B.1, depicts the relationship between the predicted and observed recidivism in the experimental sample. It is clear

that our predicted recidivism index is predictive of observed recidivism and the correlation is close to one. To obtain more power, we aggregated observed and predicted recidivism from multiple time horizons of 6, 12, 18, 24, 30, 36, 42, and 48 months from randomization.

Figure B.1: The Correlation Between Observed and Predicted Recidivism in the Experimental Sample



(a) Any rearrest



(b) Felony rearrest