

The Criminal Justice Administrative Records System: A Next-Generation Research Data Platform

Keith Finlay¹, Michael Mueller-Smith^{2,*}, and Jordan Papp³

¹U.S. Census Bureau, Washington, DC, USA

²University of Michigan, Department of Economics, Ann Arbor, MI, USA

³University of Michigan, Institute for Social Research, Ann Arbor, MI, USA

*corresponding author: mgms@umich.edu

ABSTRACT

The Criminal Justice Administrative Records System (CJARS), a joint project of the U.S. Census Bureau and the University of Michigan, is a nationally integrated data infrastructure project designed to transform research and policy-making on the United States criminal justice system. At the University of Michigan, CJARS collects longitudinal electronic records from criminal justice agencies and harmonizes these records to track a criminal episode across all stages of the system. At the U.S. Census Bureau, harmonized criminal justice records can be linked anonymously at the person-level with extensive social, demographic, and economic information from national survey and administrative records.

Introduction

In the United States, the social cost of crime is immense. In addition to the substantial costs of crime to victims, involvement in the criminal justice system not only has significant impacts on people accused or convicted of criminal offenses but also on their families and communities. Yet there is no unified data infrastructure for measuring the U.S. criminal justice system, evaluating its policies, or understanding the population that interacts with it. The lack of data infrastructure reflects the highly decentralized structure of the criminal justice system as data are held across thousands of disparate jurisdictions. There are important national data programs that cover criminal offenses or justice processes such as the National Incident-Based Reporting System (NIBRS) and National Corrections Reporting Program (NCRP), but these programs do not allow us to understand how justice system processes are connected.

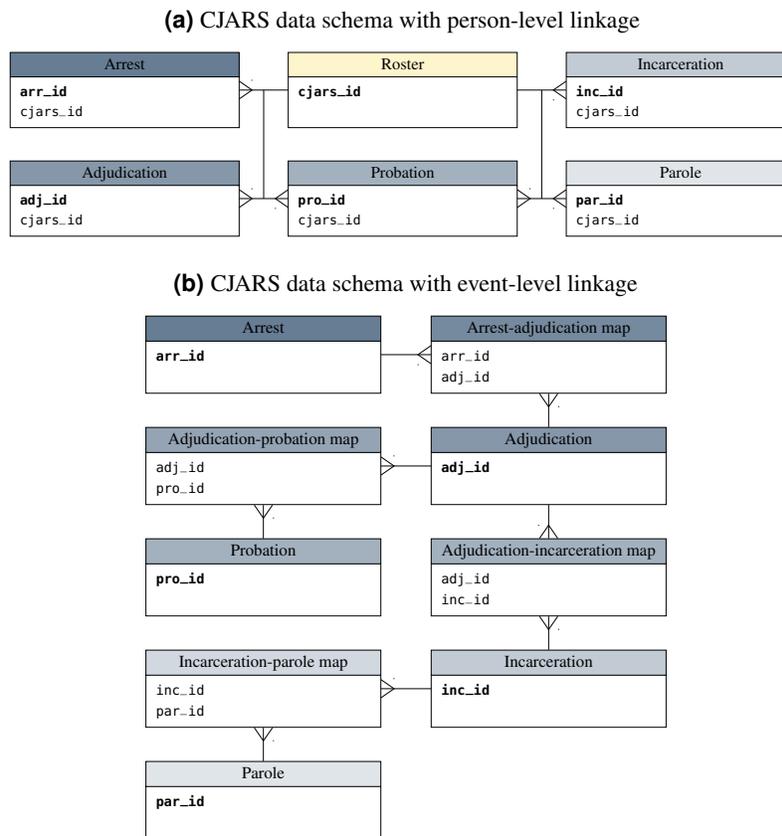
Comparable measures of criminal justice system performance require the integration of data from across the country. To identify effective policy levers, we must be able to connect criminal justice processes from arrest to sanction. Criminal justice data must be linked with other socioeconomic data to conduct dynamic benefit-cost analyses of spillover effects of criminal justice involvement. All of these types of linkages are rare, and integration with non-criminal justice data, such as labor market outcomes, is especially rare¹. To address these shortcomings, the Criminal Justice Administrative Records System (CJARS) has undertaken the task of creating an integrated criminal justice data infrastructure that can be linked at the individual level across domains of the justice system and to non-justice system outcomes². The ultimate goal of CJARS is nationwide coverage of the major types of events that occur in the justice system (i.e., arrests, criminal court case filings, and terms of probation, incarceration, and parole).

CJARS is a data collection effort and dissemination platform founded in 2016 that aims to modernize research and statistical reporting on the U.S. criminal justice system. The data infrastructure is a partnership between the University of Michigan and U.S. Census Bureau. CJARS collects, harmonizes, and integrates administrative data across five primary domains of the U.S. justice system: arrest, adjudication, incarceration, probation, and parole. The CJARS relational database schema parallels that organizational structure. Currently, CJARS includes over 2 billion lines of raw data that identify 178 million criminal justice events involving 36 million individuals across 23 states. The depth of historical data coverage varies by jurisdiction, but many states include series that extend back over 4 decades.

The U.S. lacks uniform rules across state and local jurisdictions on the privacy afforded to justice-involved individuals and what criminal justice contact is deemed public information³. Likewise, there is substantial heterogeneity in the development of data access mechanisms for researchers across the country. Lacking authority to compel data provision, CJARS relies on multiple strategies for opportunistic data acquisition. These include data use agreements, public records requests, web scraping, bulk data downloads, and data donations.

The CJARS team harmonizes disparate data sources into a standard national data schema to facilitate data integration across jurisdictions and domains of the system. Entries in each of the five CJARS procedural databases reflect events relevant to

Figure 1. CJARS relational linkage structure of roster and procedural data files



Notes: Variables shown in **bold** represent unique identifiers (primary keys) within a given dataset. Crow's feet show one-to-many table links. Researchers must first deduplicate the `cjars_id` in the Roster table before doing the person-level linkage shown in Panel A.

42 the corresponding stage of the justice system: the arrest database is measured at the arresting charge level; the adjudication
 43 database is measured at the charge level; the probation, incarceration, and parole databases are measured at the level of terms of
 44 probation, incarceration, and parole, respectively. More information about the data schema and variables available in each
 45 relational database are available in the CJARS data documentation².

46 Figure 1 provides a visual representation of the organizational structure of the CJARS data infrastructure. Panel A of
 47 Figure 1 shows the relational database structure where there is one roster file and five procedural databases (one database for
 48 each major procedural domains of the justice system). The roster file contains a unique, anonymous identifier, the `cjars_id`,
 49 which identifies individuals within the data system. Each of the five procedural databases also contains the `cjars_id` to facilitate
 50 linkage across phases of the justice system, as well as an event identifier in each database (e.g., `arr_id`) that uniquely identifies
 51 events that occur within a given domain of the justice system. Panel B of Figure 1 demonstrates how the procedural databases
 52 can be linked in conjunction with associative tables to reconstruct the sequence of events that are related to a single criminal
 53 episode. More information about the data schema and variables available in each relational database can be found in the CJARS
 54 data documentation².

55 CJARS data are available to qualified researchers on approved projects through the Federal Statistical Research Data Centers
 56 (FSRDCs), a network of 32 secure physical locations where CJARS can be linked to other anonymized survey and administrative
 57 data held in the Census Bureau's Data Linkage Infrastructure. Researchers must apply and have projects approved before
 58 access, but the FSRDCs provide a proven, secure mode of data distribution that can safeguard the privacy and confidentiality
 59 of the extensive micro data contained in CJARS. It also provides the additional benefit of enabling individual-level linkage
 60 with a range of other socioeconomic data held in the FSRDC network, including self-reported demographic characteristics,
 61 evolving family composition and place of residence, employment and earnings behavior, take-up of public benefit programs,
 62 and mortality. For example, CJARS can be linked at the individual-level to responses from decennial census or American
 63 Community Survey (ACS) responses. The scope of the CJARS data holdings and other data available through the FSRDC
 64 network provide extensive opportunities for researchers to conduct novel research that previously was not possible without such

65 a data infrastructure.

66 CJARS is continually expanding its geographic and procedural coverage through its data collection efforts. Continued data
67 collection will provide broader coverage and thus improved capacity to support research. New vintages of the CJARS data
68 infrastructure are made available in the FSRDC network on an approximately bi-annual basis.

69 **Results**

70 Due to variation in data collection methods and the numerous creative solutions required to coherently process the data collected
71 by CJARS, there is a fundamental need to benchmark the data infrastructure against other available data series to both validate
72 the strengths of CJARS and to highlight its potential weaknesses to interested researchers. The degree of data integration
73 available through the CJARS micro data is unprecedented, creating challenges for ideal technical validation testing. Our
74 validation efforts focus mainly on reproducing available aggregate statistical series published by the Bureau of Justice Statistics
75 (BJS), with the implicit assumption that success in matching aggregate information bolsters the validity of the underlying micro
76 data contained in CJARS as well.

77 We evaluate CJARS against the following federal statistical series: Uniform Crime Report (UCR)⁴, State Court Processing
78 Statistics Series (SCPS)⁵, National Prisoners Statistics Program (NPS)⁶, National Corrections Reporting Program (NCRP)⁷,
79 Annual Probation Survey⁸, and Annual Parole Survey⁹. These programs share with CJARS a common set of demographic and
80 criminal justice measures in a common set of state and local jurisdictions with considerable historical data to assess data quality
81 over time. Our focus is to benchmark CJARS data at the state-level rather than to aggregate across all CJARS states whenever
82 possible. This reflects the decentralization of the U.S. criminal justice system, and allows us to better assess our data collection
83 and harmonization practices. Our evaluation of arrest data against the UCR can only be made at local and county levels, so we
84 omit it here. Please refer to Papp & Mueller-Smith [10] for these findings.

85 Our analyses focus on reproducing caseload count and flow estimates (e.g., yearly entries into prison as measured in the
86 NPS), as well as caseload characteristics and outcomes (e.g., demographic characteristics of defendants in SCPS data) in
87 CJARS-covered jurisdictions. CJARS-based estimates that closely corroborate existing federal estimates provide important
88 evidence on the quality and accuracy of our nascent data infrastructure endeavor and the population-level, linkable micro-data
89 from which the CJARS-based estimates are constructed.

90 **Comparing CJARS Adjudication to SCPS**

91 The U.S. lacks a comprehensive statistical reporting program on the criminal court system. The closest option we have
92 for validating CJARS adjudication records is the SCPS program, an occasionally produced statistical series that documents
93 characteristics of felony defendants from large urban counties. A number of common caseload composition and case processing
94 metrics can be calculated using both CJARS and SCPS, creating the opportunity to benchmark the CJARS adjudication data for
95 the subset of records that overlaps with the definition of the scope of SCPS. This still provides useful information on gauging
96 the quality of the algorithms applied to all of our data. Examples include average age of defendants, defendant gender and
97 race/ethnicity, disposition type, time between disposition and sentencing, probation and incarceration sentence length, and
98 offense type.

99 Figure 2a provides a scatter plot where comparable SCPS and CJARS statistics (e.g., average age of felony defendants) are
100 plotted onto the y- and x-axes, respectively. Individual statistics are plotted for each of the 1996, 1998, 2000, 2002, 2004, 2006,
101 and 2009 waves of SCPS with corresponding CJARS statistics built from the same corresponding time frames. The color/shape
102 of a marker in the scatter plot represents a specific outcome and are repeated across the survey waves. The expectation is that
103 the plotted points will cluster around the reference line which has a slope equal to one. Clustering around the line indicates that
104 the statistics generated using CJARS and SCPS are comparable.

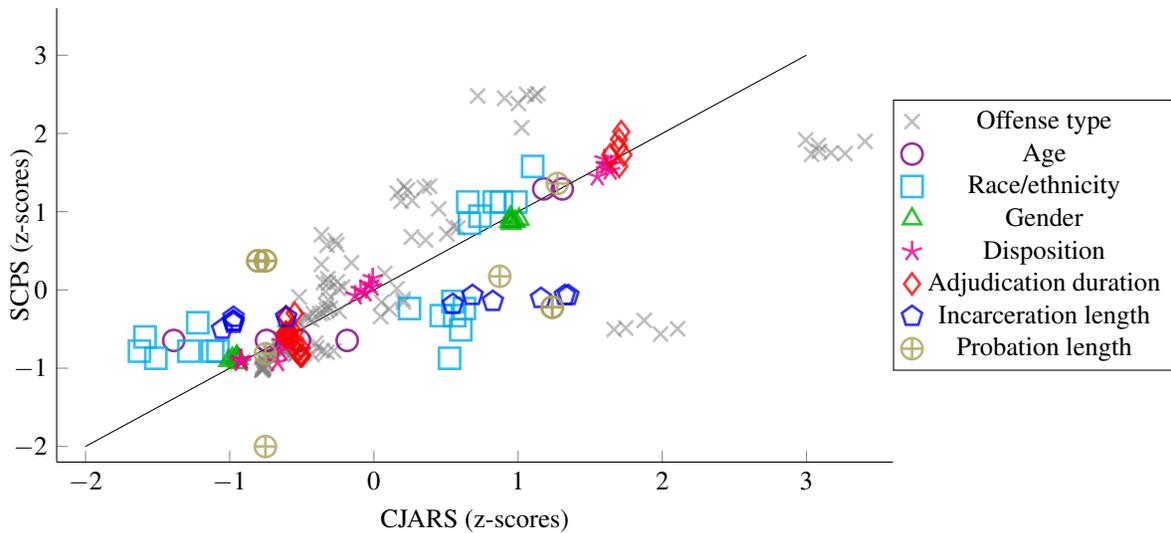
105 Figure 2b shows formal tests of the difference between means calculated using SCPS microdata and means calculated using
106 CJARS case data. The means have been calculated over jurisdiction-year cells, which represents a total of 609 jurisdiction-year
107 combinations covered by each wave of the SCPS data. In order to account for differential temporal and jurisdictional coverage,
108 the tests first residualize for event-year by jurisdiction fixed effects meaning the tests evaluate the statistical similarity of
109 caseload means between SCPS and CJARS within overlapping jurisdiction-year cells. Since jurisdictions are not resampled
110 within CJARS, we cluster our standard errors at the jurisdiction level to account for their repeated observation over time. On all
111 but one metric (Offense type: violent), we cannot reject the null hypothesis that SCPS and CJARS-based measures are equal,
112 and in this case it is only marginally significant (p-value = 0.059).

113 **Comparing CJARS Incarceration to the NPS and NCRP**

114 CJARS incarceration records can be compared to similar data from the NPS and the NCRP. All three data sources contain
115 information that can be used to estimate annual prison entry counts, exit counts, and populations as well as incarceration rates.
116 For succinctness, we present comparisons here only for annual entry counts.

Figure 2. CJARS and SCPS-derived statistics for felony defendants in large urban counties

(a) Plotting standardized CJARS and SCPS caseload statistics, by jurisdiction-year



(b) Testing differences of means between SCPS and CJARS, by jurisdiction-year

	Age	White	Black	Hispanic	Other race	Male	Female
<i>Panel A: Defendant demographics</i>							
Mean difference (SCPS-CJARS)	0.083 (0.263) [0.752]	-0.081 (0.098) [0.408]	-0.035 (0.044) [0.426]	0.106 (0.095) [0.265]	-0.0003 (0.003) [0.933]	-0.018 (0.017) [0.316]	0.018 (0.017) [0.316]
<i>Panel B: Offense type</i>	Violent	Property	Drug	Public order			
Mean difference (SCPS-CJARS)	0.096 (0.050) [0.059]	-0.129 (0.087) [0.141]	0.067 (0.046) [0.151]	-0.034 (0.065) [0.608]			
<i>Panel C: Disposition outcomes</i>	Days between disposition and sentencing	Disposition: diversion	Disposition: dismissal	Disposition: conviction	Sentence: incarceration (months)	Sentence: probation: (months)	
Mean difference (SCPS-CJARS)	-9.55 (12.12) [0.434]	0.038 (0.042) [0.369]	0.006 (0.108) [0.954]	-0.026 (0.110) [0.817]	-1.62 (6.59) [0.807]	-4.40 (5.70) [0.443]	

Source: Calculations from CJARS data held by the University of Michigan and not protected by 13 USC §9a¹⁰.

Panel A notes: This figure compares average caseload characteristics of CJARS adjudication microdata filed in May of the years 1996, 1998, 2000, 2002, 2004, 2006, and 2009 from jurisdictions representative of the 75 largest counties to the average caseload characteristics reported in all waves of the SCPS series. SCPS samples felony filings in May of the reported years. Comparisons are made on offense type, defendant gender and race/ethnicity, average defendant age, incarceration length, probation length, disposition type, and length of time between disposition and sentencing. Differences between CJARS and SCPS generated statistics were transformed into z-scores (where $z = (x_i - \mu) / \sigma$) to improve readability.

Panel B notes: Each regression is estimated on 609 jurisdiction-year observations, weighted by caseload size. Testing the mean difference between SCPS and CJARS caseload statistics is evaluated using a binary indicator for whether the statistic originated from SCPS (=1) or CJARS (=0), controlling for a fully saturated set of jurisdiction by event-year fixed effects. Standard errors, clustered at the jurisdiction level, are displayed in parentheses. P-values are shown in brackets.

Figure 3 provides a comparison of annual entry counts as reported in the NPS and the NCRP, and from calculations using CJARS. A separate graph is given for each state for which CJARS has historical data holdings. In each graph, the purple line represents CJARS, blue the NPS, and green the NCRP. CJARS closely aligns with either the NPS, NCRP, or both in nearly every graph. For example, annual entry counts align well in Pennsylvania between all three data sources. Conversely, CJARS aligns better with either the NPS or NCRP but not both in, for example, Washington and North Carolina. A similar set of exercises that compare annual exit counts, year-end populations, and incarceration rates showing substantively similar findings can be found in Papp & Mueller-Smith [10].

Additionally, we calculated the absolute average annual percent difference in the CJARS, NPS, and NCRP prison entry counts on a state-by-state basis. Across all years and states, CJARS entry counts differ from NPS and NCRP entry counts by an absolute average difference of 15.8% and 11.9%, respectively. In comparison, NPS and NCRP entry counts differ from each other by an absolute average difference of 16.1%—a larger discrepancy than CJARS has with either series.

Comparing CJARS Probation and Parole to the Annual Probation and Parole Surveys

The probation information in CJARS provides information that can be compared to similar data from the Annual Probation Survey. Both CJARS and the Annual Probation Survey can be used to estimate yearly entry and exit counts as well as yearly probationer populations and rates. Comparisons here focus on entry counts for succinctness.

Figure 4 shows a comparison between probation entry counts observed in CJARS as compared to the Annual Probation Survey for each state where CJARS has historical data holdings. The graphs show substantial alignment in North Carolina. There also appears to be good alignment in Michigan, but there is considerable instability from year-to-year entry counts in the Annual Probation Survey leading to large increases and decreases. In comparison, the CJARS data from Michigan provide much more stable counts from year-to-year. The graph for Texas in Figure 4 shows similarities when coverage in the CJARS data begins (early 2000s). However, a gap forms over time in which more entries are observed in the CJARS data. A similar set of exercises that compare annual exit counts, year-end populations, and probationer rates showing substantively similar findings can be found in Papp & Mueller-Smith [10].

Additionally, we calculated the average yearly percent difference on a state-by-state basis between CJARS and the Annual Probation Survey in terms of probation entry counts. Then, we calculated the absolute mean yearly difference across all CJARS-covered states and years to quantify the average difference in entry counts between CJARS and the Annual Probation Survey. Comparing CJARS to the Annual Probation Survey shows an absolute average difference of 14.4%.

The parole information in CJARS provides information that can be compared against the same types of information gathered as part of the Annual Parole Survey. Both sources of parole data can be used to estimate yearly entry and exit counts as well as yearly parolee populations and rates. Comparisons here focus on entry counts for succinctness.

Figure 5 shows a comparison between parole entry counts observed in CJARS as compared to the Annual Parole Survey for each state where CJARS has historical data holdings. As can be seen in this figure, entry counts in CJARS and the Annual Parole Survey line up exceptionally well in almost all states. The one state where there is a slight difference is Nebraska where the counts of events in CJARS are slightly lower than those reported in the Annual Parole Survey. However, the difference is consistent across years and so the trends of changes in entry counts over time align between CJARS and the Annual Parole Survey. A similar set of exercises that compare annual exit counts, year-end populations, and parolee rates showing substantively similar findings can be found in Papp & Mueller-Smith [10].

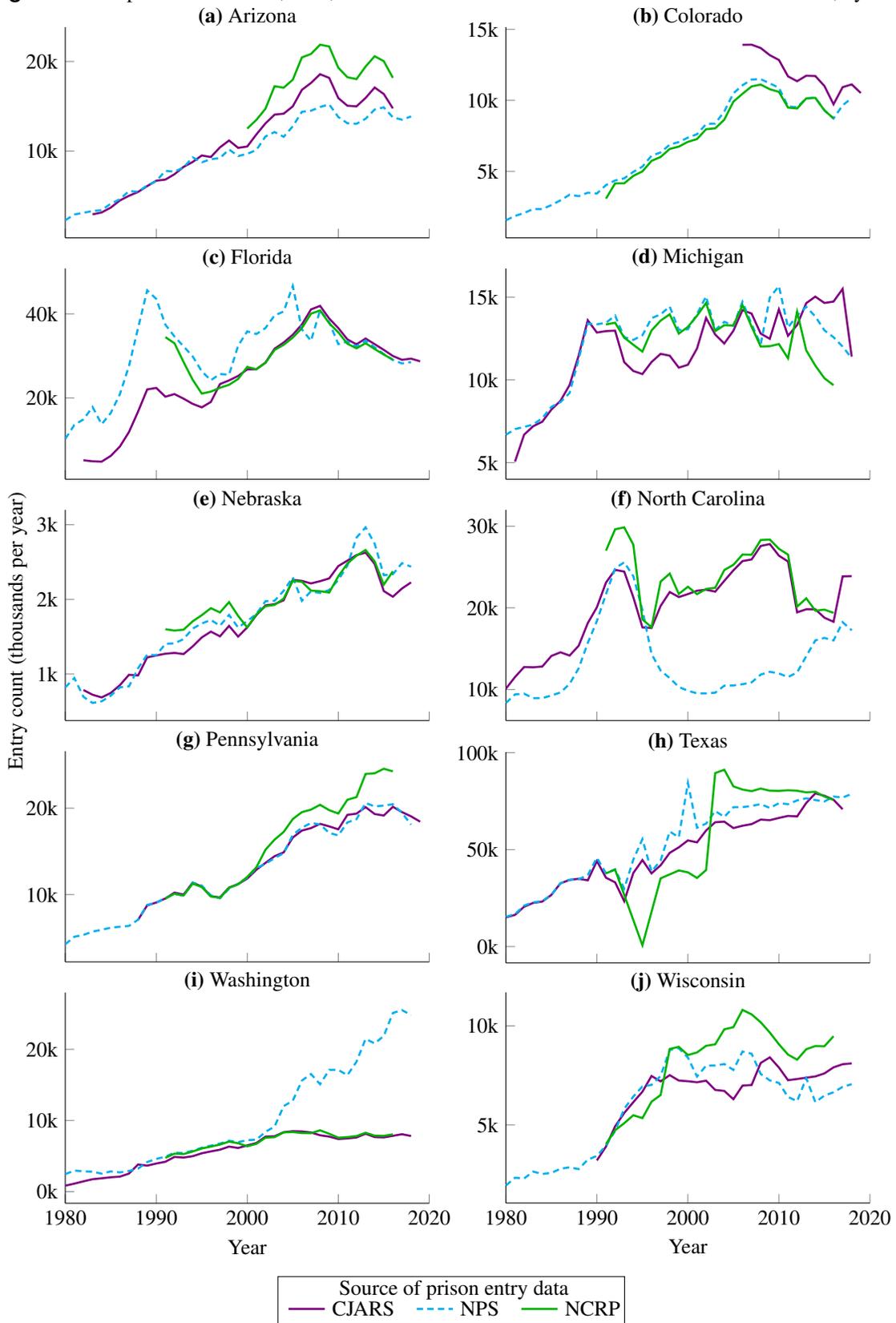
Finally, we calculated the absolute average yearly percent difference on a state-by-state basis between CJARS and the Annual Parole Survey in terms of parole entry counts. Then, we calculated the mean difference across all CJARS-covered states and years to quantify the average difference in entry counts between CJARS and the Annual Parole Survey. Comparing CJARS to the Annual Parole Survey shows an absolute average difference of 15.7%.

Discussion

Given that CJARS data quality is validated through closely replicating extant federal statistical series, the research data platform represents a transformative resource to advance knowledge on the determinants of criminal activity and the individual and community impacts of the U.S. criminal justice system. No existing repository available to researchers (1) is composed of records that cover criminal justice agencies of all types and from all geographies; (2) measures criminal justice events from arrest through sanction at the person level for the entire population; and (3) can be linked with extensive information about the socioeconomic characteristics and outcomes of justice-involved individuals.

CJARS data holdings continue to grow, but the data platform does not yet cover the entire country, raising questions about the appropriateness for population level statistics. Consequently, it is worth considering whether a CJARS dataset without complete national coverage is representative of the criminal justice system more broadly. In Figure 6, we compare CJARS-covered states to non-CJARS covered states along three dimensions: average violent crime rates between 2000 and 2018⁴, average property crime rates between 2000 and 2018⁴, and average imprisonment rates between 2000 and 2018¹¹.

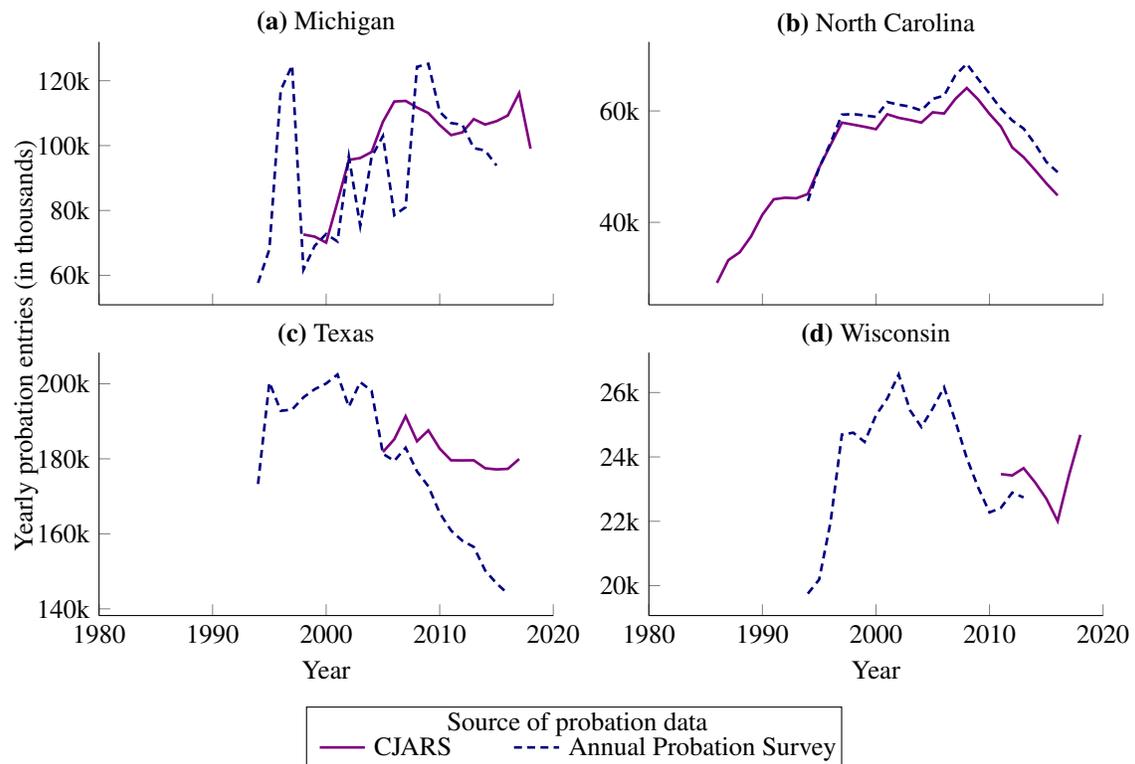
Figure 3. Comparison of CJARS, NPS, and NCRP-based estimates of annual incarceration entries, by state



Source: Calculations from CJARS data held by the University of Michigan and not protected by 13 USC §9a¹⁰.

Notes: These graphs compare the yearly counts of individuals entering prison in a given state. Each graph covers the period of time between 1980 and 2020. However, historical coverage of data varies across states for CJARS, the NPS, and the NCRP. Yearly entry counts were estimated by creating a count of individuals that entered prison in a given year using CJARS and NCRP microdata, respectively. Yearly entry counts are reported using aggregate estimates from the NPS.

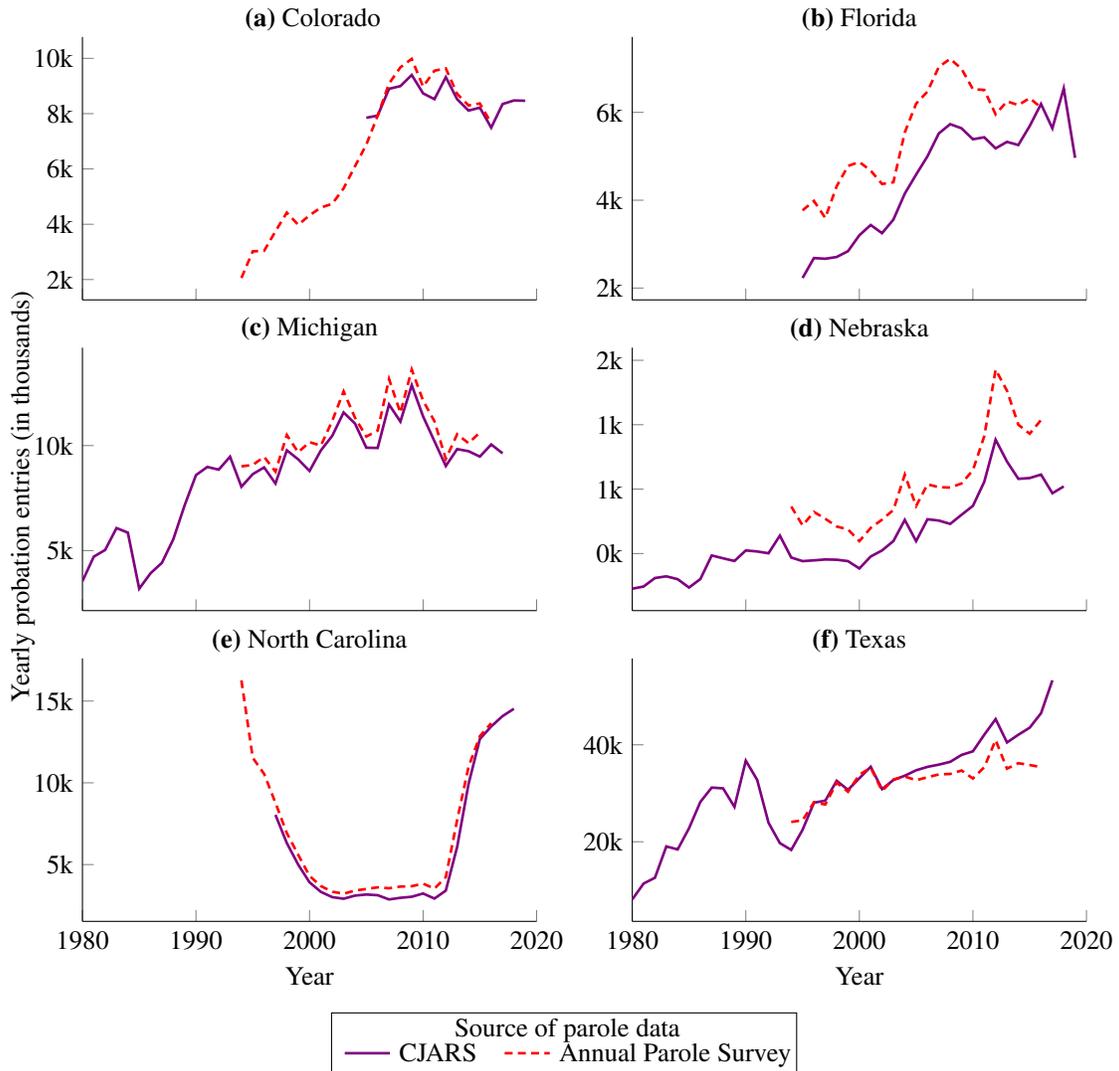
Figure 4. Comparison of CJARS and Annual Probation Survey-based estimates of annual probation entries, by state



Source: Calculations from CJARS data held by the University of Michigan and not protected by 13 USC §9a¹⁰.

Notes: These graphs compare the yearly counts of individuals beginning probation terms in a given state. Each graph covers the period of time between 1980 and 2020. However, historical coverage of data varies across states for CJARS and the Annual Probation Survey. Yearly entry counts were estimated by creating a count of individuals that began probation in a given year using CJARS microdata. Yearly entry counts are reported using aggregate estimates from the Annual Probation Survey.

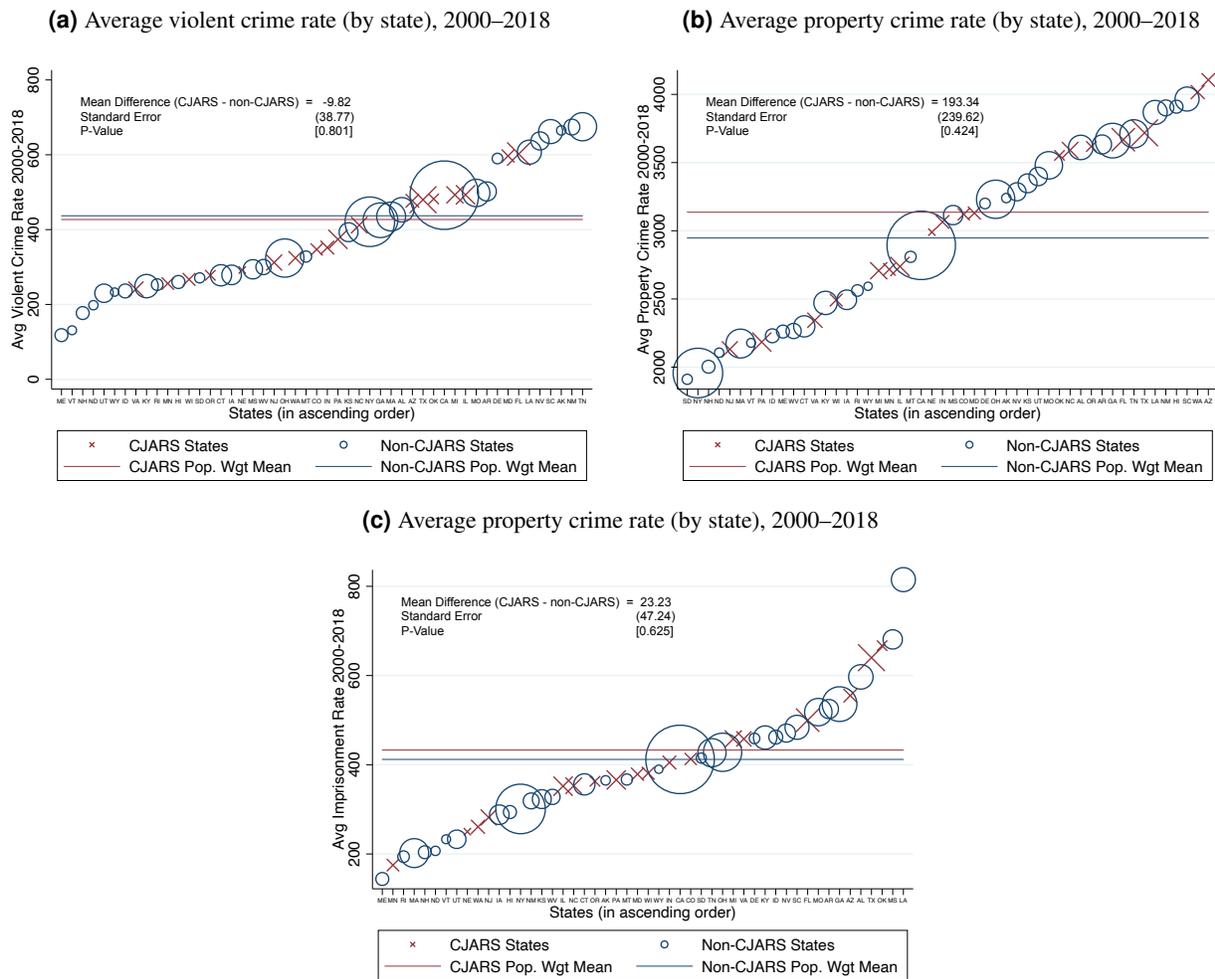
Figure 5. Comparison of CJARS and Annual Parole Survey-based estimates of annual parole entries, by state



Source: Calculations from CJARS data held by the University of Michigan and not protected by 13 USC §9a¹⁰.

Notes: These graphs compare the yearly counts of individuals beginning parole terms in a given state. Each graph covers the period of time between 1980 and 2020. However, historical coverage of data varies across states for CJARS and the Annual Parole Survey. Yearly entry counts were estimated by creating a count of individuals that began parole in a given year using CJARS microdata. Yearly entry counts are reported using aggregate estimates from the Annual Parole Survey.

Figure 6. Representativeness of CJARS data coverage for UCR-measured crime rates and BJS-measured imprisonment rate



Source: Calculations from publicly available BJS¹¹ and UCR⁴ annual reports. Calculations from CJARS data held by the University of Michigan and not protected by 13 USC §9a¹⁰.

Notes: Mean differences between CJARS and non-CJARS covered states are weighted according to each state’s average population size between 2000 and 2018. Standard errors are shown in parentheses and p-values from a two-sided test evaluating the null hypothesis of zero difference are in brackets.

170 From this exercise, we observe two key facts that bolster the case that CJARS can provide useful national estimates while
 171 continuing to grow toward national coverage. First, for each comparison series, the difference in the weighted means between
 172 CJARS and non-CJARS covered states is not statistically significantly different from zero. In fact, the differences in means are
 173 modest compared to weighted means observed in non-CJARS states: -2.2% for violent crime rates, 6.6% for property crime
 174 rates, and 5.6% for imprisonment rates. Second, we can also see in the figures that CJARS-covered states are represented
 175 throughout the distributions of each of the comparison measures, which suggests that CJARS data holdings can provide a
 176 representative perspective on a broad range of criminal justice processes.

177 **Methods**

178 **Ethics**

179 Justice-involved people are a vulnerable population, and a core principle of the CJARS project is that we acquire, store, and
 180 analyze criminal justice data securely and ethically so that the identities and characteristics of the individuals in the CJARS data
 181 are kept confidential. The CJARS data collection and repository was reviewed and approved by the University of Michigan
 182 Health Sciences and Behavioral Sciences Institutional Review Board (approval number REP00000094); the review included
 183 additional oversight from a prisoner advocate per federal regulations. The project received a waiver for informed consent

184 because it is built from existing electronic records maintained by government entities and involves no direct contact with or
185 interventions applied to any of human subjects. Any further research activity involving CJARS, including research conducted
186 with anonymized CJARS data through the secure FSRDC network, requires separate Institutional Review Board approval. The
187 validation results described in this paper were separately approved as research activity (approval number HUM00208278).

188 **Data Collection**

189 The CJARS project collects data from police departments, sheriff offices, prosecutors, criminal courts, departments of
190 corrections, and state criminal history repositories. The primary and preferred method of data collection is signed data-use
191 agreements with agencies. These legal agreements delineate the responsibilities of the University of Michigan and the allowable
192 uses of acquired data. Each of these agreements authorizes the University of Michigan to transfer data securely to the Census
193 Bureau for linkage-based research and statistical work.

194 The CJARS team at the University of Michigan also collects data by submitting public records requests, sometimes referred
195 to as Freedom of Information Act (FOIA) requests, in states where statutes require or allow agencies to make criminal history
196 information available to the public. In a third approach, the University of Michigan harvests data that are publicly available
197 online using web scrapers or bulk downloads. All web scraping is carried out under a set of ethical scraping policies to ensure
198 that collection complies with agency website robots.txt and Terms of Use.

199 The goal of CJARS is to integrate data from federal, state, and local agencies. To maximize the growth of data coverage
200 while using project resources efficiently, CJARS prioritizes acquisitions from agencies that manage statewide data systems
201 including departments of corrections, state court administrative offices, and state criminal history repositories. CJARS does
202 acquire local agency data where the costs of doing so are low, such as where web scraping or public records requests are
203 possible. As a benchmark, CJARS aims to spend no more than \$0.01 per acquired row of data, which typically shifts our focus
204 to larger jurisdictions where increasing returns to scale reduce per-observation acquisition costs.

205 **Data Processing**

206 One of the major barriers to research on the criminal justice system is a lack of data integration across agencies. The CJARS
207 team implements a systematic set of procedures to process and link the data it collects into a single, integrated data platform.
208 Figure 7 provides a visual depiction of this process which will be used to describe CJARS data processing in the following
209 sections.

210 **Data Processing at the University of Michigan**

211 Data collected from data providers are initially stored on secure data servers at the University of Michigan. Original data are
212 cleaned and harmonized by the CJARS team at the University of Michigan. Cleaning and harmonization is an extensive process
213 that involves transforming data received in its raw form from data providers into a format that fits the CJARS data schema.
214 CJARS employs several strategies to conduct data processing and harmonization.

215 CJARS data processing at the University of Michigan is broken out into a sequence of six steps (see Figure 8a). First, after
216 data has been collected from a data provider, native data formats undergo *localization* to apply a common database format to
217 each individual dataset. Second, the *standardization* stage extracts and harmonizes personally identifying information (PII),
218 and imputes gender and race/ethnicity information where needed. PII variables are then used as inputs in *entity resolution*
219 to generate a unique, person-level identifier that tracks involvement in the justice system across jurisdictions, over time, and
220 through the various procedural domains of the justice system. Our approach to entity resolution leverages a biometrically
221 trained probabilistic matching model¹².

222 The entity resolution process identifies and assigns each individual a unique `cjars_id` (see Figure 8b). These identifiers are
223 added to the cleaned data that have been stripped of all PII variables and are transferred to the anonymized partition of the
224 CJARS secure data servers, where the records undergo further processing. The separation between the PII and anonymized
225 partitions aims to restrict access to PII variables to CJARS staff with an operational need, adding additional privacy and
226 confidentiality protections within our organization.

227 With `cjars_ids` attached, the data next proceeds through *harmonization*, which brings each individual dataset into the
228 common national schematic adopted in the CJARS data platform. The purpose of harmonization is to align disparate source
229 files to reduce barriers for multi-jurisdictional research. This is accomplished through both populating a uniform set of variables
230 from each source file and ensuring coded values follow a consistent standard. One example of the latter is offense classification
231 where we have to translate over 4 million unique text descriptions into a unified set of offense codes. This specific task is
232 accomplished through a machine learning model that CJARS has developed in partnership with Measures for Justice, known as
233 the Text-based Offense Classification (TOC) tool¹³.

234 The final two data processing steps at the University of Michigan involve *event deduplication* and *episode resolution*.
235 Because we receive data from agencies with overlapping data coverage (e.g., statewide repositories and local criminal courts)
236 as well as repeated extracts over time with evolving information on local caseloads, the de-duplication stage is critical to ensure

Figure 7. CJARS stakeholders, data exchange, record harmonization, and product development

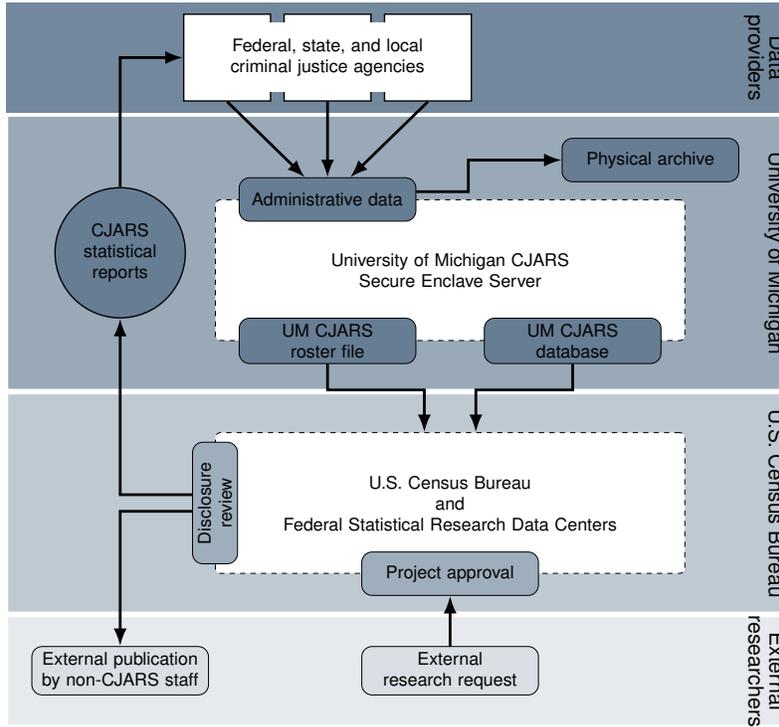
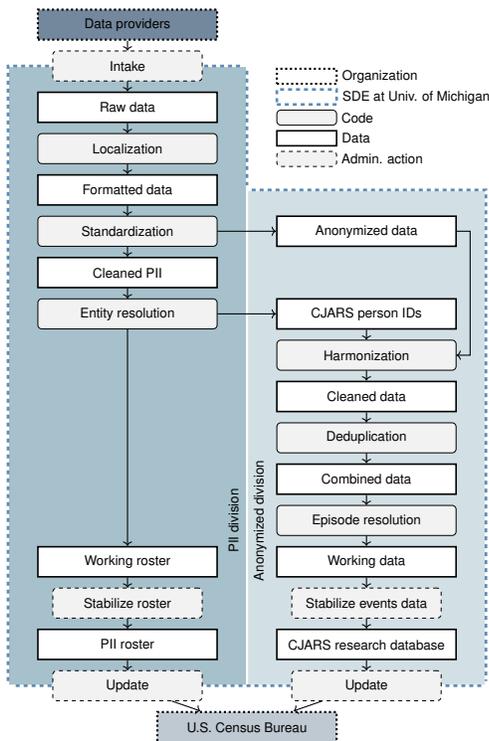
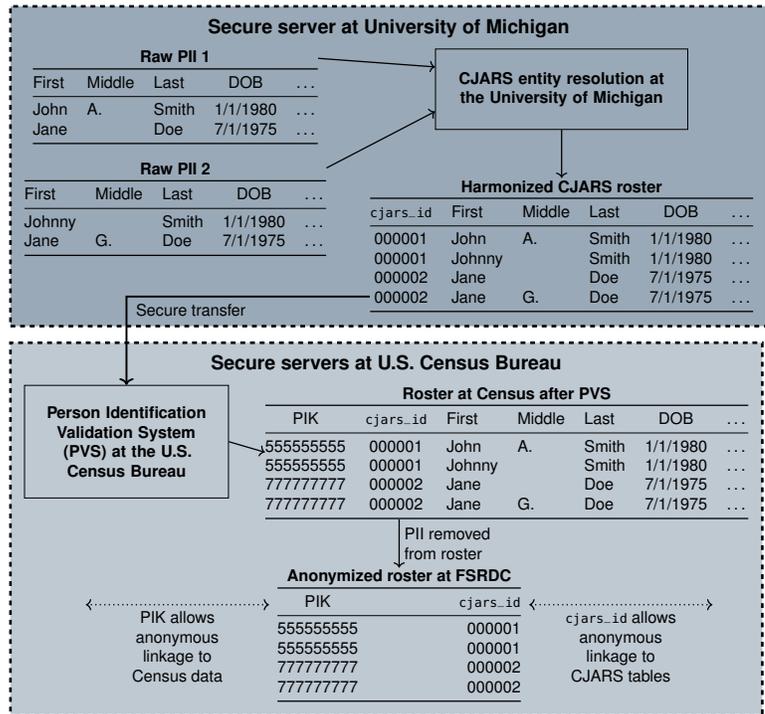


Figure 8. CJARS data processing steps

(a) Data processing sequence at the University of Michigan



(b) Entity resolution and record linkage



237 that we are not over-counting the number of distinct points of contact that individuals have with the justice system. Lastly,
238 episode resolution generates crosswalks that connect the procedural stages of the justice system to each other using contextual
239 information like event timing, offense types, and sentencing outcomes.

240 **Data Integration at the Census Bureau**

241 After data processing at the University of Michigan is complete, the data are securely transferred to the U.S. Census Bureau for
242 integration into the FSRDC system. This begins with the CJARS roster file being processed through the Person Identification
243 Validation System (PVS)¹⁴, a probabilistic record linkage system that generates a crosswalk between the `cjars_id` and the
244 Protected Identification Key (PIK). PIKs uniquely and anonymously identify individuals in the U.S. within the FSRDC system,
245 and allow researchers to link de-identified CJARS and non-CJARS survey and administrative records at the individual-level.
246 Research occurs within a secure computing environment that is available at Census Bureau headquarters and in the FSRDCs
247 across the U.S.

248 **Data Availability**

249 CJARS data may be accessed through the FSRDC network, which is composed of 32 secure physical locations where qualified
250 researchers on approved projects can link anonymized data from the Census Bureau's Data Linkage Infrastructure. Many
251 FSRDC projects are currently approved for virtual access.

252 All five CJARS procedural databases and the `cjars_id-to-PIK` crosswalk file (also known as the anonymized roster file)
253 are available to qualified researchers with approved Census Bureau projects in its FSRDC network. Data are only accessible
254 through the FSRDC network to ensure the privacy and security of the sensitive records contained in CJARS. Anonymized data
255 are stored in SAS format, although the FSRDC network has a wide range of statistical software available to support researchers
256 working with their preferred tools.

257 The CJARS team has produced extensive documentation and support tools to assist data users. The CJARS data documenta-
258 tion provides information on project scope, data collection methods, data coverage, data access and security, data processing
259 and harmonization, linkage techniques, and data providers. The documentation explains the data schema and includes a variable
260 codebook with variable-level descriptive statistics and data notes describing unique aspects of specific data acquisitions that
261 data users should consider.¹

262 Prospective researchers can apply to use CJARS through the FSRDC project application process, which begins by speaking
263 with an FSRDC administrator. The CJARS team has produced a proposal guide that walks researchers through the process,
264 identifies who can apply, describes the criteria used to approve projects, and documents data-use limitations.²

265 **Code Availability**

266 CJARS project code that does not contain sensitive information is made available at <https://github.com/umcjars>. This
267 includes a shell of our data production process, which can be found at [https://github.com/umcjars/cjars_production_](https://github.com/umcjars/cjars_production_code_shell)
268 [code_shell](https://github.com/umcjars/cjars_production_code_shell). In addition, code for the validation exercise presented in this article can be found at [https://github.com/](https://github.com/umcjars/cjars_bjs_validation)
269 [umcjars/cjars_bjs_validation](https://github.com/umcjars/cjars_bjs_validation). In addition, code used for other research projects involving CJARS data conducted at the
270 Census Bureau through the FSRDC system is available upon request via email at erd.cjars@census.gov. This code must be
271 cleared first by the Census Bureau to avoid disclosure of any potentially sensitive information.

272 **References**

- 273 1. Plecas, D., McCormick, A. V., Levine, J., Neal, P. & Cohen, I. M. Evidence-Based Solution to Information Sharing
274 between Law Enforcement Agencies. *Policing: An International Journal of Police Strategies and Management*. <https://doi.org/10.1108/13639511111106641> (2011).
- 275 2. Finlay, K. & Mueller-Smith, M. *Criminal Justice Administrative Records System (CJARS)* Ann Arbor, MI: University of
276 Michigan. 2021. <https://cjars.isr.umich.edu/data-documentation-download/>.
- 277 3. Jacobs, J. & Crepet, T. The Expanding Scope, Use, and Availability of Criminal Records. *NYU Journal of Legislation*
278 *and Public Policy* **11**, 177. [https://nyujlpp.org/wp-content/uploads/2012/10/Jacobs-Crepet-The-Expanding-](https://nyujlpp.org/wp-content/uploads/2012/10/Jacobs-Crepet-The-Expanding-Scope-Use-and-Availability-of-Criminal-Records.pdf)
279 [Scope-Use-and-Availability-of-Criminal-Records.pdf](https://nyujlpp.org/wp-content/uploads/2012/10/Jacobs-Crepet-The-Expanding-Scope-Use-and-Availability-of-Criminal-Records.pdf) (2007).
- 280 4. Federal Bureau of Investigation. "Crime in the U.S." Data Series 2000-2018. [https://www.fbi.gov/services/cjis/](https://www.fbi.gov/services/cjis/ucr/publications)
281 [ucr/publications](https://www.fbi.gov/services/cjis/ucr/publications).

¹The data documentation can be downloaded from <https://cjars.isr.umich.edu/data-documentation-download>.

²The list of FSRDC locations can be found at <https://www.census.gov/about/adrm/fsrdc>. The CJARS proposal guide can be found at <https://cjars.isr.umich.edu/proposal-guide-download>.

- 283 5. Bureau of Justice Statistics. *State Court Processing Statistics, 1990–2009* 2014. <https://doi.org/10.3886/ICPSR02038.v5> (2021).
284
- 285 6. Bureau of Justice Statistics. *National Prisoner Statistics* 1978-2018. <https://doi.org/10.3886/ICPSR37639.v1> (2021).
286
- 287 7. Bureau of Justice Statistics. *National Corrections Reporting Program* 1978-2018. <https://doi.org/10.3886/ICPSR37973.v1> (2021).
288
- 289 8. Bureau of Justice Statistics. *Annual Probation Survey* 1994-2018. <https://www.icpsr.umich.edu/web/NACJD/series/327> (2021).
290
- 291 9. Bureau of Justice Statistics. *Annual Parole Survey* 1994-2018. <https://www.icpsr.umich.edu/web/NACJD/series/328> (2021).
292
- 293 10. Papp, J. & Mueller-Smith, M. *Benchmarking the Criminal Justice Administrative Records System's Data Infrastructure* Working paper. 2021. <https://cjars.isr.umich.edu/benchmarking-report-download>.
294
- 295 11. Bureau of Justice Statistics. *"Prisoners in" Publication Series* 2000-2018. https://bjs.ojp.gov/library/publications/list?series_filter=Prisoners.
296
- 297 12. Gross, M. & Mueller-Smith, M. *Modernizing Person-Level Entity Resolution with Biometrically Linked Records* Working
298 paper. 2020. https://sites.lsa.umich.edu/mgms/wp-content/uploads/sites/283/2020/12/entity_resolution_20201203.pdf.
299
- 300 13. Choi, J., Kilmer, D., Mueller-Smith, M. & Taheri, S. *Hierarchical Approaches to Text-based Offense Classification*
301 Working paper. 2021. http://sites.lsa.umich.edu/mgms/wp-content/uploads/sites/283/2022/01/CJARS_MFJ_offense_classification_20220119.pdf.
302
- 303 14. Wagner, D., Lane, M., et al. *The Person Identification Validation System (PVS): Applying the Center for Administrative
304 Records Research and Applications' (CARRA) Record Linkage Software* Census Bureau CARRA Working Paper 2014-01.
2014. <https://www.census.gov/library/working-papers/2014/adrm/carra-wp-2014-01.html>.

305 Acknowledgements

306 The project has been supported by National Science Foundation grant SES-1925563, as well as grants from the Laura and
307 John Arnold Foundation, the Bill and Melinda Gates Foundation, and the Robert Wood Johnson Foundation. The University of
308 Michigan has supported the CJARS project through the following programs: Michigan Institute for Teaching and Research
309 (MITRE), Populations Studies Center (PSC), and Poverty Solutions. Finlay's work at the Census Bureau has been supported by
310 funding from the Evidence-Based Policymaking Commission Act (P.L. 114-140) and the Foundations for Evidence-Based
311 Policymaking Act (P.L. 115-435). Any conclusions expressed are those of the authors and do not necessarily represent the
312 views of the U.S. Census Bureau.

313 Author Contributions

314 Mueller-Smith and Finlay co-founded the CJARS project. Mueller-Smith directs the CJARS project and oversees all aspects
315 of data collection, processing, harmonization, and research at the University of Michigan. Finlay co-directs the project and
316 oversees data integration at the U.S. Census Bureau. Papp conducted the technical validation analyses. All authors assisted in
317 the drafting and editing of the manuscript.

318 Competing Interests

319 The authors declare no competing interests.