

# BIG DATA: A Tool for Fighting Discrimination and Empowering Groups

---









## PREFACE

In May 2014, the Executive Office of the President concluded its 90-day study of Big Data and privacy and released a report entitled *Big Data: Seizing Opportunities, Preserving Values*. The report highlighted certain positive uses of Big Data, such as identifying health risks at an early stage, creating efficiencies in energy distribution, and uncovering fraud through predictive analysis. However, it also concluded that Big Data analytics' could facilitate discrimination in housing, credit, employment, health, education, and a range of other markets. These potential benefits and drawbacks underscore the need to better understand how Big Data will shape our lives in years to come.

Recognizing the 50<sup>th</sup> anniversary of the Civil Rights Act and the challenges to fighting discrimination in the 21<sup>st</sup> century, the case studies included in this report show how businesses, governments, and civil society organizations are leveraging Big Data (and other data sets) to protect and empower vulnerable groups, including by providing access to job markets, uncovering discriminatory practices, and creating new tools to improve education and assist those in need. While by no means an exhaustive list of Big Data's potential to uncover and fight discrimination, we offer these examples to show how Big Data already is redefining efforts to ensure equal opportunity for all.

We would like to thank the Anti-Defamation League for its partnership in preparing this report, Jared Bomberg and Julian Flamant at Hogan Lovells US LLP for providing essential research and drafting support, and members of the FPF Advisory Board for reviewing drafts and providing guidance.

We hope these examples will contribute to discussions about Big Data's impact on discrimination.

Jules Polonetsky  
Executive Director and Co-Chair  
Future of Privacy Forum

Christopher Wolf  
Founder and Co-Chair  
Future of Privacy Forum







# TABLE OF CONTENTS

I. Seeing Beyond Bias to Provide New Opportunities .....	1
Case Study 1: Workplace Diversity (Entelo) .....	1
Case Study 2: Smarter Hiring (Google) .....	2
Case Study 3: Allocation of Public Works (Cedar Grove Institute) .....	3
Case Study 4: Demographics of Health (State of New York) .....	4
II. Transparency is a Necessary Disinfectant .....	5
Case Study 5: Auditing Anti-Semitism (Anti-Defamation League) .....	5
Case Study 6: <i>McClesky v. Kemp</i> .....	6
Case Study 7: Discrimination Complaint Data (EEOC) .....	7
Case Study 8: <i>United States v. Sterling</i> .....	8
III. Long-Term Problems Require New Solutions .....	9
Case Study 9: Education for All (NSBA) .....	9
Case Study 10: Tracking Migratory Patterns (United Nations) .....	10
Case Study 11: Economic Development and Equality (OECD) .....	11
Case Study 12: Finding Missing and Exploited Children (Palantir/NCMEC) .....	12
Case Study 13: Human Trafficking (Palantir/Polaris Project) .....	13



## I. SEEING BEYOND BIAS TO PROVIDE NEW OPPORTUNITIES

### Case Study 1: Workplace Diversity (Entelo)

*Entelo Diversity*, a candidate recruiting platform launched in April 2014, is improving workplace diversity by empowering recruiters to search for job candidates from within underrepresented segments of the population. Using a proprietary algorithm, this workplace diversity tool sifts through publicly available data—pulled from social media platforms—to match recruiters with candidates who hold necessary qualifications, but also meet particular diversity requirements. The tool can filter candidates based on gender, race, and military history in five categories: Female, African American, Asian, Hispanic, and Veteran.

#### Example of an Entelo Search

The screenshot displays the Entelo search interface. On the left, a search bar contains the keyword 'PHP'. Below it, filters are applied for 'San Francisco, CA 50mi' and 'Email available'. The 'Diversity' filter is expanded, showing 'Female' selected. The main area shows '3,807 candidates found' and a list of three candidate profiles:

- Eve Killaby**: SAN FRANCISCO, CALIFORNIA. Back-end Engineer at Bloodhound - about 1 year. Senior Web Engineer at DocuSign, Inc. - 9 months. State University of New York at Buffalo, Bachelors of Science, Computer Science. Skills: mysql, node.js, software engineering, web development, PHP, web design, entrepreneurship, 25 more.
- Lia Napolitano**: SAN FRANCISCO BAY AREA. Siri Interaction Designer at Apple Inc. - over 1 year. User Experience Designer at Apple Inc. - almost 2 years. Wellesley College, BA, Media Arts and Sciences. Skills: user interface design, user interface, mac, user-centered design, PHP, flash, iOS, 31 more.
- Julia Kryzstofiak-Szopa**: SAN FRANCISCO BAY AREA, US. Founder at Wellfitting.com - about 1 year. Captain Ivanova at Blackbox Accelerator, LLC - about 1 year. Katholieke Universiteit Leuven, Philosophy, AI. Skills: git, html 5, web development, social media, PHP, user experience, Django, 51 more.

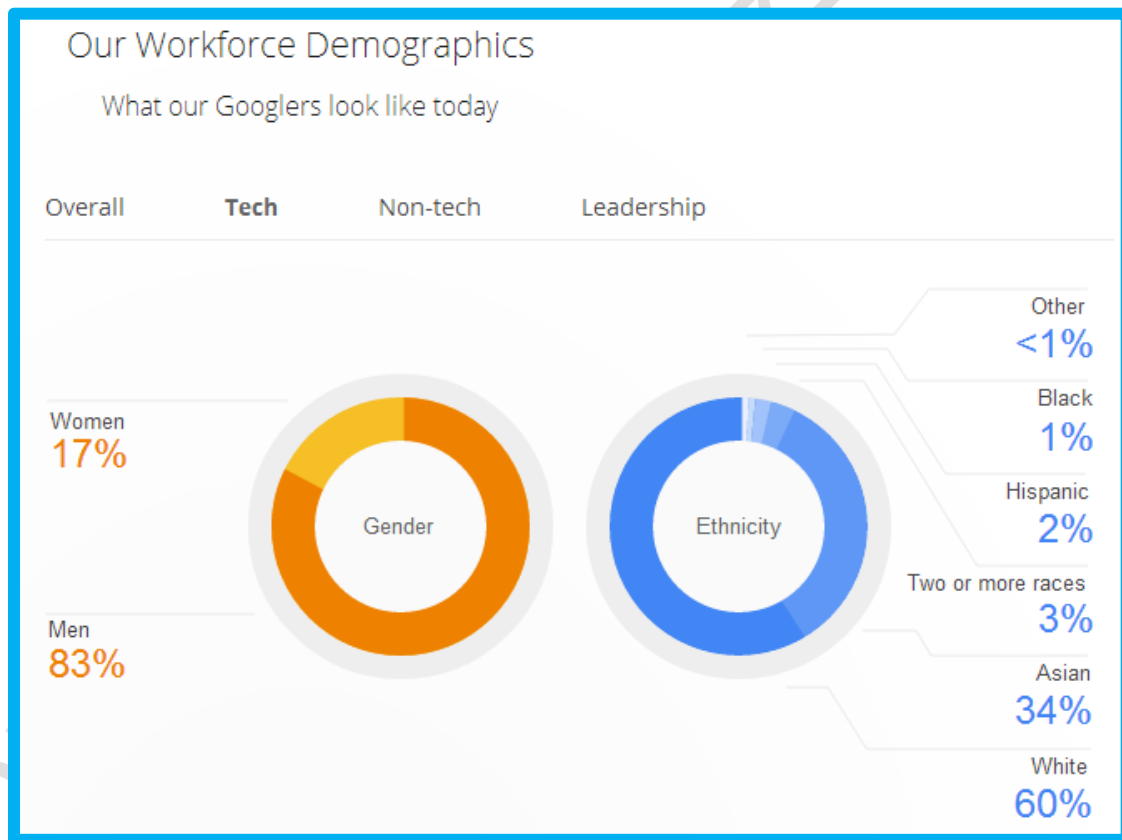
Source: [Entelo](http://Entelo)



## Case Study 2: Smarter Hiring (Google)

A challenge for the technology industry is ensuring diversity in the workplace. Twitter has recently reported that 90% of its global “tech” employees are male and Google admits “[it’s] not where [it] wants to be when it comes to diversity,” with only 1% black and 17% women among its tech workforce.

Recognizing the value of a diverse workforce, [Google](#) is leveraging its data analytics capabilities to help change those numbers. Through analytics, the company identified issues with its traditional hiring process, which favors certain groups and may provide little insight into a candidate’s “Googliness.” Based on this analysis, the company recently implemented a new interview process focused on structured behavioral questions, rather than “brainteasers” and academic G.P.A’s. Google also now uses employee panels to evaluate candidates through several rounds of review. The result of Google’s intensive analytical study and the adoption of new hiring methods is that the company is better able to evaluate candidates on an individualized basis, where questions like race and gender do not hold sway.



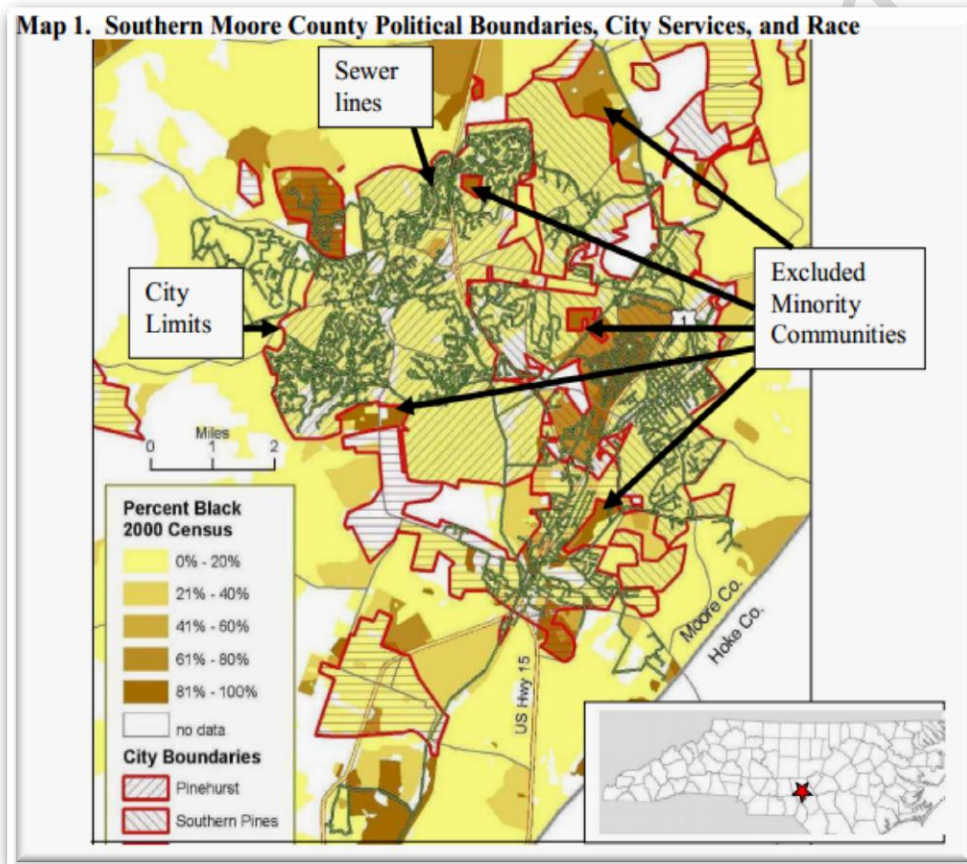
Source: [Google Diversity](#)



### Case Study 3: Allocation of Public Works (Cedar Grove Institute)

The Cedar Grove Institute for Sustainable Communities is a non-profit organization that leverages open data to explore disparities in allocation of geographic boundaries and public works of various communities. Cedar Grove uses a combination of demographic analysis, contextual investigation, housing and economic analysis, and geographic information systems to explore potentially discriminatory implications of public policy decisions. In a 2005 study, [\*Segregation in the Modern South: A Case Study of Southern Moore County\*](#), Cedar Grove Institute combined census data and other publicly available surveys and demographic information to explore the impact on community development of land annexation policies of Moore County North Carolina.

#### Map Showing Excluded Minority Communities



Source: [Cedar Grove Institute for Sustainable Communities](#)



## Case Study 4: Demographics of Health (State of New York)

In 2011, the Institute of Medicine, the health arm of the National Academy of Sciences, [reported](#) that lesbian, gay, bisexual, and transgender (LGBT) individuals have unique health experiences and needs, but as a nation, we do not know exactly what these needs are. The IOM also reported that clinicians and researchers are faced with incomplete information regarding the health status of LGBT individuals and that current research has not adequately examined subpopulations, particularly racial and ethnic groups and peoples' health needs based on age.

In response, the State of New York has [launched](#) a coordinated, multi-agency effort to strengthen data collection regarding LGBT individuals in New York. The campaign will rely on data collected on a self-reporting basis by the New York's Department of Health, Department of Corrections and Community Supervision, Office for the Aging, Office of Mental Health, Office of Alcohol and Substance Abuse Services, Office of Temporary and Disability Assistance, Office of Children and Family Services, and Office for People with Developmental Disabilities. The data collected will be shared among the eight agencies to create a comprehensive method for identifying the needs of the lesbian, gay, bisexual, and transgender community. It is hoped that stronger data sets will empower the State and others to create more tailored approaches to reduce health disparities impacting the LGBT community.

### News Article

Wednesday, July 23, 2014

### GOVERNOR CUOMO ANNOUNCES MULTI-AGENCY STATE EFFORT TO ADDRESS LGBT DISPARITIES

*New York becomes first state in the nation with coordinated statewide strategy to improve LGBT data collection*

Governor Andrew M. Cuomo announced that New York State is undertaking a coordinated, multi-agency effort to strengthen data collection for lesbian, gay, bi-sexual and transgender (LGBT) New Yorkers. Outlined in the first report by the State's Interagency LGBT Task Force, this statewide effort to include sexual orientation and gender identity information in data collections will allow the state to better tailor services to meet LGBT needs, ultimately improving the health and lives of thousands of New Yorkers. This effort makes New York the first state in the nation to employ a coordinated strategy to develop its data collection procedures for the LGBT community.

"New York State has a long history of advancing progressive ideals, and today we are continuing to lead the nation by identifying new ways to improve services and better meet the needs of the LGBT community," Governor Cuomo said. "By being more inclusive with how state agencies monitor the demographics of those they serve, we can address health and financial disparities, safety concerns, and a myriad of other issues that impact LGBT New Yorkers. This is another step forward for an important community in New York, and our administration will continue standing up for all New Yorkers, regardless of their sexual orientation or gender identity."

The Institute of Medicine in its March 2011 report, *The Health of Lesbian, Gay, Bisexual, and Transgender People: Building a Foundation for Better Understanding*, emphasized the need for collection of gender identity

Source: [State of New York](#)



## II. TRANSPARENCY IS A NECESSARY DISINFECTANT

### Case Study 5: Auditing Anti-Semitism (Anti-Defamation League)

Anti-Semitism is a persisting global problem. Each year the Anti-Defamation League publishes an *Audit of Anti-Semitic Incidents*, which reports on both criminal and non-criminal acts of harassment and intimidation, including distribution of hate propaganda, threats and slurs. In 2013 alone, there were 751 anti-Semitic incidents across the United States. The annual report combines data on anti-Semitic assaults, vandalism, and harassment events as collected by 27 regional offices across the U.S. While the data is broken down to detail each event and show granular state-wide trends, it is also aggregated to provide a comprehensive analysis of anti-Semitism on a nationwide basis. In 2013, total anti-Semitic incidents declined by 19 percent across the U.S., while “violent” incidents increased by 17 percent.

In recent years, the ADL Audit has placed anti-Jewish incidents into one of three categories: assault, vandalism or harassment. In 2013 the numbers broke down as follows:

- *Assaults*: 31 incidents reported in 2013, compared with 17 in 2012;
- *Vandalism*: 315 incidents in 2013, compared with 440 in 2012;
- *Harassment, threats and events*: 405 incidents in 2013, compared with 470 in 2012.

#### State-By-State Totals

While the number of anti-Semitic incidents in the U.S. decreased in 2013, ADL recorded increases in some states when compared with 2012 figures, such as in Connecticut, Michigan, Ohio and Texas. Additionally, continuing a consistent trend for many years, the states with the highest totals were those with large Jewish populations. Once again, New York and California, while each experiencing declines compared to the previous year, topped the list:

- *New York State*, with 203 incidents in 2013, down from 245 in 2012;
- *California*, with 143 incidents, down from 185;
- *New Jersey*, with 78 incidents, down from 173;
- *Florida*, with 68 incidents, down from 88;
- *Massachusetts*, with 46 incidents, up from 38;
- *Pennsylvania*, with 43 incidents, up from 37.

Source: [Anti-Defamation League](#)



## Case Study 6: *McClesky v. Kemp*

McCleskey, an African American man, was sentenced to death after being convicted of armed robbery and the murder of a white police officer. In a writ of habeas corpus, McCleskey argued that the Georgia capital sentencing process was administered in a racially discriminatory manner in violation of the Eighth and Fourteenth Amendments. In support of the claim, McCleskey offered a statistical study (the Baldus study) to show disparities in the imposition of the death sentence in Georgia based on the murder victim's race and the defendant's race. The study was based on over 2,000 murder cases that occurred in Georgia during the 1970's, and involved data relating to the victim's race, the defendant's race, and the various combinations of such persons' races. The study found a consistent pattern of discrimination in the use of the death penalty against defendants who were charged with killing white victims compared to defendants who were charged with killing African American victims.

While the court ultimately found against McCleskey, the case has been described as a turning point in the debate over the death penalty in the United States. The Baldus study has been replicated in numerous jurisdictions with similar findings. Race is now a powerful issue in debates over the death penalty because of studies like the Baldus study, which show that race can affect death penalty decisions.



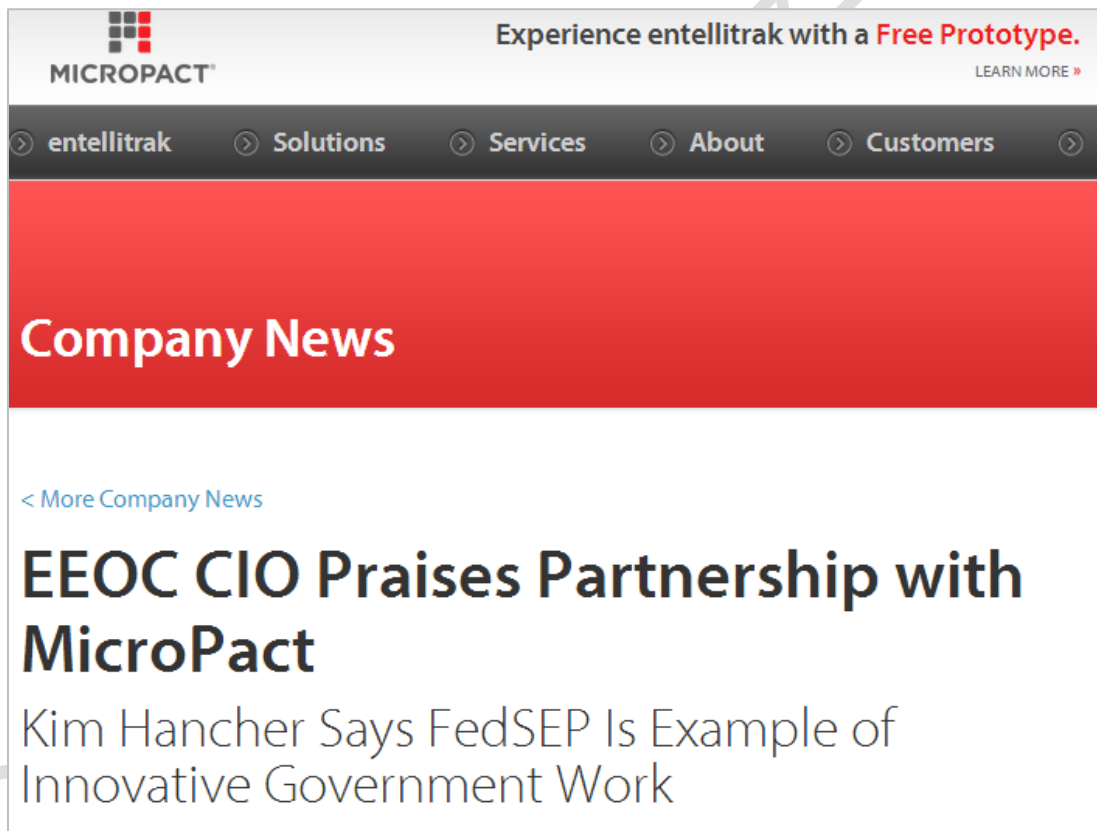
Source: *McClesky v. Kemp*, 481 U.S. 279 (1987).



## Case Study 7: Discrimination Complaint Data (EEOC)

The Equal Employment Opportunity Commission (EEOC) is responsible for enforcing federal laws that make it illegal to discriminate against a job applicant or an employee because of the person's race, color, religion, sex (including pregnancy), national origin, age (40 or older), disability, or genetic information.

In March 2013 the EEOC unveiled FedSEP, an electronic portal through which more than 325 federal agencies interact with the EEOC relative to their workforce and complaint data. This gateway provides each agency's staff with a single point of access to EEOC data collection systems and provides a new tool to collect and analyze government-agency data on workplace discrimination charges. By streamlining the submission process for so many documents and aggregating data from many sources, FedSEP allows EEO professionals greater ability to spot trends and uncover discrimination across the federal government.



Source: [Entellitrak](#)



## Case Study 8: *United States v. Sterling*

In 2006, Donald Sterling, his wife, and their family-trust real-estate company, *Beverly Hills Properties*, were accused of engaging in discriminatory practices in violation of the Fair Housing Act and Title VIII of the Civil Rights Act of 1968. The United States alleged, *inter alia*, that the defendants refused to rent portions of their 27-building development in the “Koreatown” neighborhood of Los Angeles to non-Koreans (e.g., *Hispanics and blacks*). A study by Dr. Shelley Lapkoff used a database of tenant information released by the defendants to show that the number of Korean tenants across the 27 buildings had increased “significantly,” from 64 percent to 83 percent, within the year following acquisition of those buildings by the defendants. Dr. Lapkoff’s report also included an analysis of census data to determine whether the defendant’s claim that the changing demographic distribution of Koreatown could explain the decreasing diversity of its tenants.

Dr. Lapkoff’s analysis found that the overall demographics of Koreatown remained relatively stable during the period in question, and that Hispanics remained the dominant race in the area. The report concluded that, absent external changes, the increase in Korean tenants was consistent with the United States’ allegation of housing discrimination. A later study by Dr. Lapkoff also used census data to show that there were no major shifts in household income or Korean households in the area that could explain the increase of Korean renters.

The studies helped lead to a settlement agreement that included a number of measures aimed at ending discriminatory renting practices in the 27 buildings owned by Beverly Hills Properties and required the defendants to pay \$2,625,000 to be disbursed among aggrieved persons and a \$100,000 civil penalty.

Table 3

Percentage of New Tenants Who Were Korean, Before and After Acquisition

Building #	Year Before Acquisition				Year After Acquisition		
	Total	Korean	% Korean		Total	Korean	% Korean
Buildings Where Less than 80 percent of New Tenants are Korean at Time of Acquisition							
101	12	0	0%		16	4	25%
112	4	0	0%		6	6	100%
93	14	1	7%		18	13	72%
96	9	1	11%		11	9	82%
102	25	3	12%		22	8	36%
105	8	2	25%		7	5	71%
108	20	5	25%		30	28	93%
85	12	4	33%		14	13	93%
106	20	7	35%		18	11	61%
98	7	4	57%		11	11	100%
88	43	26	60%		41	34	83%
97	8	5	63%		10	8	80%
82	29	20	69%		27	25	93%
95	33	24	73%		44	38	86%
104	19	15	79%		28	25	89%
Subtotal	263	117	44%		303	238	79%
Buildings Where More than 90 percent of New Tenants are Korean at Time of Acquisition							
87	40	36	90%		45	43	96%
89	12	11	92%		13	13	100%
83	11	11	100%		9	9	100%
84	18	18	100%		17	14	82%
90	8	8	100%		19	17	89%
91	9	9	100%		4	4	100%
94	5	5	100%		10	10	100%
103	14	14	100%		18	16	89%
107	14	14	100%		10	9	90%
109	14	14	100%		20	18	90%
110	12	12	100%		7	6	86%
Total	157	152	97%		172	159	92%

Building # 86 Had No Rentals the Year Before Acquisition for Comparison

Source: *United States v. Sterling*, No. 2:06CV04885, (C.D. Cal. Nov. 12, 2009).



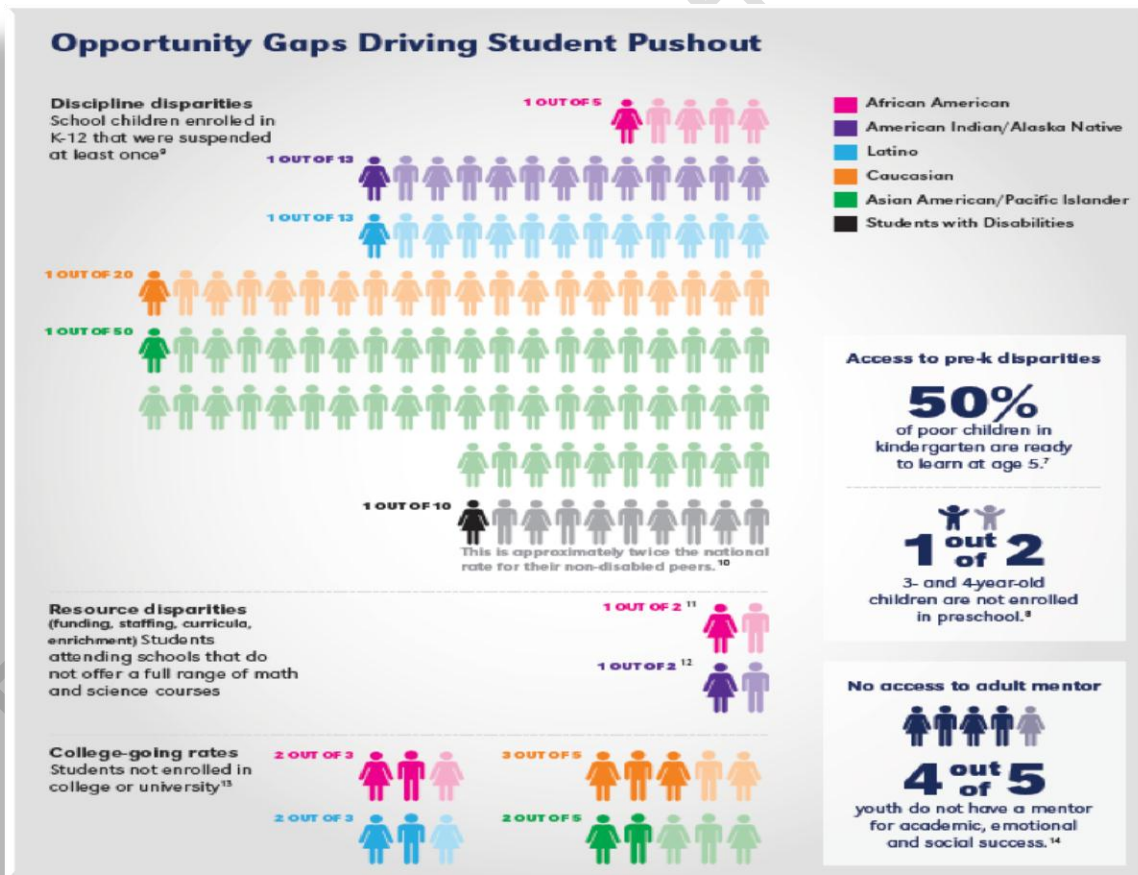
### III. LONG-TERM PROBLEMS REQUIRE NEW SOLUTIONS

#### Case Study 9: Education for All (NSBA)

A recent report by National School Boards Association (NSBA) offers novel policy solutions for increasing education rates in America. The report, *Partnerships, not Pushouts*, combines census data with data collected by various organizations to identify factors—known as “pushouts”—that may be responsible for driving young people away from education. Pushout factors can be more common among different segments of the population. For example, school suspensions—considered a major “pushout” factor, affect one out of five African American students and only one out twenty Caucasian students, which may partly explain the large discrepancy between graduation rates of those two groups.

To increase education levels among American youth, the NSBA proposes a variety student-centered “Personal Opportunity Plans” (POPs). To be effective, POPs are tailored to meet the needs of students on an individualized basis, addressing the pushout factor(s) most threatening to a particular student’s academic success.

#### Discipline Disparities



Source: [National School Board Association](#)

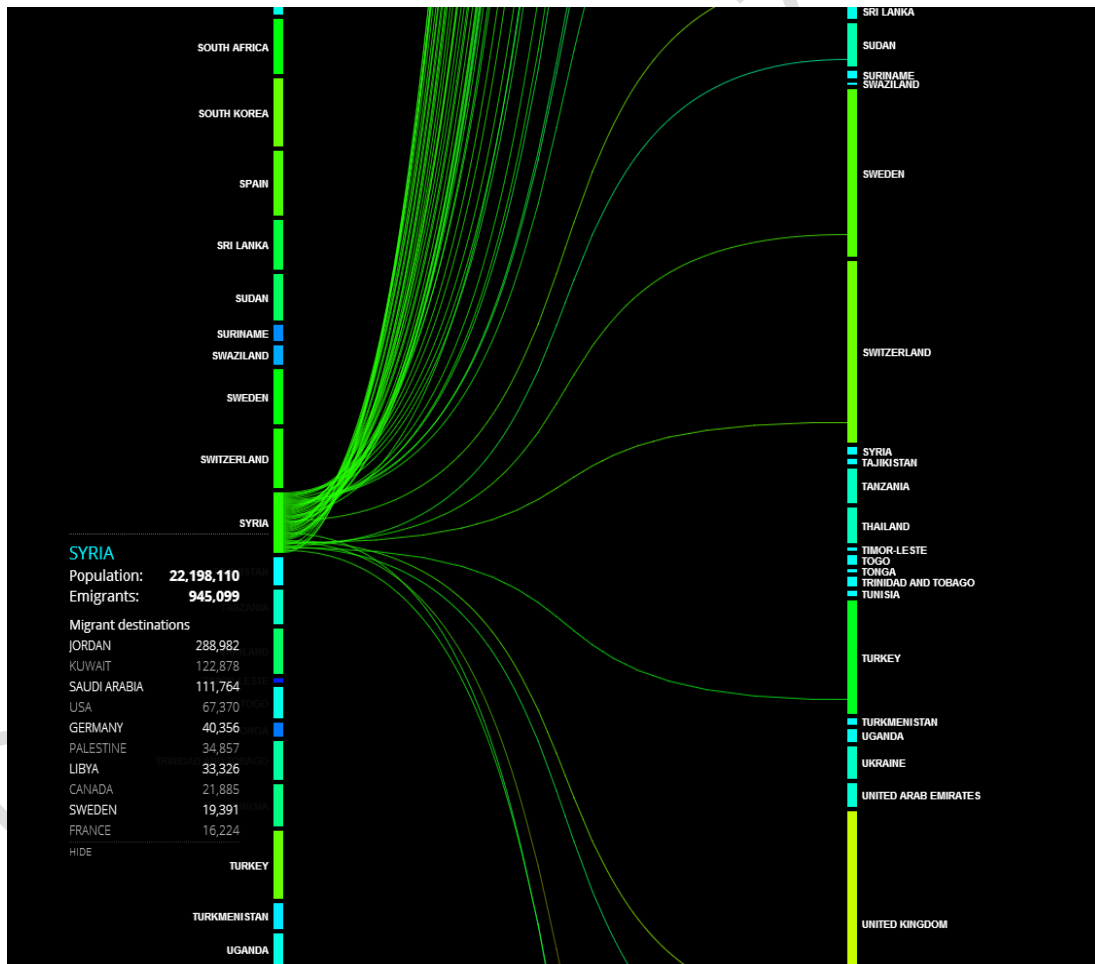


## Case Study 10: Tracking Migratory Patterns (United Nations)

The [United Nations has highlighted the social benefits of tracking migratory patterns](#) of diverse peoples. For example, tracking the movement of displaced populations can empower humanitarian groups to provide better aid to those populations. As a new project under the *UN Global Pulse* banner, the organization is exploring new ways to track displaced populations using Big Data. The organization cites “significant shortcomings” with traditional methods of migratory benchmarking such as censuses, demographic and thematic surveys and administrative registers, which quickly become outdated.

In its review, Global Pulse highlights a number of studies that rely on Big Data collected from social media sites or open data initiatives to draw important conclusions about population movements. In the example below, PeopleMovin, repurposes “open” migration, refugee and asylum, and world population data to create an interactive tool allowing users to quickly identify international movement patterns and identify where relief efforts are most valuable.

### Migration Patterns from Syria:



Source: [http://peplemov.in/#f\\_SY](http://peplemov.in/#f_SY)



## Case Study 11: Economic Development and Equality (OECD)

Since 2009, the Organization for Economic Co-Operation and Development has offered the publicly available *Gender, Institutions and Development Database*. The database compiles gender-discrimination data from 160 countries to provide researchers and policymakers with an analysis of 60 detailed variables, ranging from factors like “Discriminatory Family Code” to “Restricted Civil Liberties,” that are likely to impact women’s engagement in society and the economy.

A defining feature of the GID-DB is that, in addition to traditional quantitative analyses, the database uses an innovative scoring system to evaluate discriminatory institutional features. For example, while traditional studies of “early marriage” analyze rates of marriage among various age groups, the GID-DB has created a scaled system that combines rates of early marriage with an analysis of legal, traditional and religious customs to provide a much deeper look at gender discrimination. An example of the scaled system is provided below:

- 0: The law on the minimum age of marriage does not discriminate against women.
- 0.5: The law on the minimum age of marriage discriminates against some women, for example through customary, traditional and religious law.
- 1: The law on the minimum age of marriage discriminates against all women or there is no law on the minimum age of marriage.

### Gender, Institutions and Development Database (GID-DB)

The screenshot displays the OECD.StatExtracts interface for the Gender, Institutions and Development Database (GID-DB). The interface includes a search bar, navigation tabs, and a table of data for various countries. The table is titled "Gender, Institutions and Development Database 2012 (GID-DB)" and shows data for the year 2012. The table includes columns for Region, Income group, Year, and various variables related to gender discrimination. The variables are grouped into categories: Discriminatory Family Code, Restricted Physical Integrity, and Social Institutions. The table shows data for countries including Afghanistan, Albania, Algeria, Angola, and Argentina.

Country	Legal Age of Marriage	Early Marriage	Parental Authority	Inheritance	Violence Against Women (Laws)	Female Genital Mutilation	Domestic Violence	Reproductive Integrity
Afghanistan	1	0.17	1	0.5	0.5	0.75	0.5	0.25
Albania	0	0.1	0	0	0.5	0.25	0.5	0.08
Algeria	0	0.02	1	0	0.5	0.5	0.75	0.25
Angola	0.5	0.36	0	0	0.5	0.25	1	0.78
Argentina	0	0.11	0	0	0	0.25	0.25	0.06

Source: [OECD, Gender, Institutions and Development Database \(GID-DB\)](#)

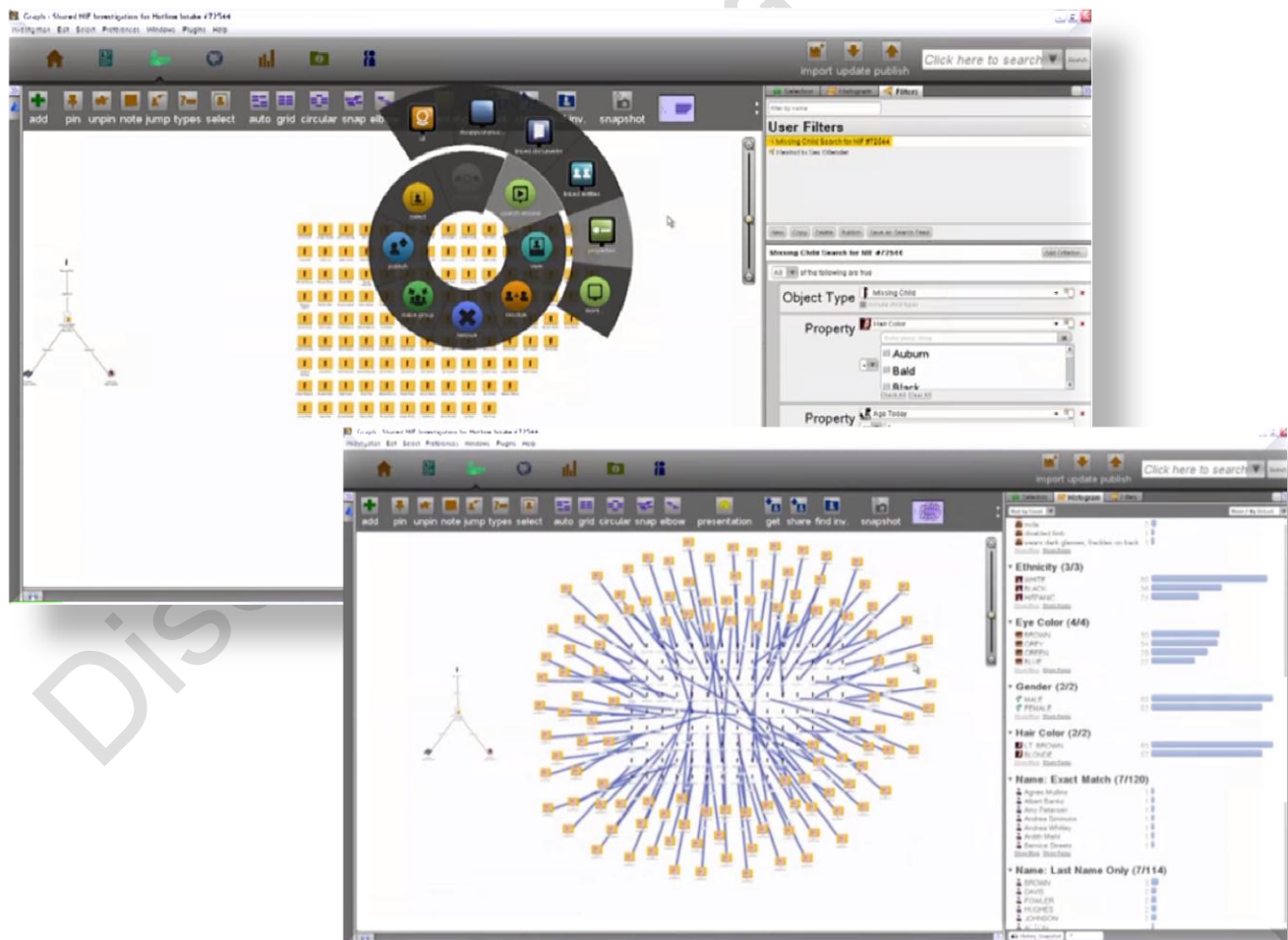


## Case Study 12: Finding Missing and Exploited Children (Palantir/NCMEC)

The National Center for Missing & Exploited Children (NCMEC) compiles a wide variety of information from law enforcement, social media, and proprietary databases. Much of this information has traditionally been stored in siloed databases, requiring analysts to manually query each database when investigating a case. The Big Data analytics tool, developed in 2010 by Palantir, empowers NCMEC analysts to query a range of databases simultaneously.

The below case illustrates how the NCMEC uses Big Data to save children:

A 17-year-old girl was reported missing and suspected of being a victim of sex trafficking. Through various searches, a NCMEC analyst was able to find multiple posts online that advertised this missing child for sex. Through information in the ads, the analyst was able to tie them to other posts from the same pimp. The analysis included over 50 advertisements, 9 different females, and a trail covering 5 states. A Link Analysis graph was created using Palantir that helped law enforcement to easily see the large scope of the ring. This insight helped law enforcement link the pimp to a multitude of other crimes and other girls that he victimized.



Source: [Palantir](#), [Using Palantir to Uncover Hidden Links In Missing And Exploited Children Cases](#)



## Case Study 13: Human Trafficking (Palantir/Polaris Project)

Human Trafficking is a global problem. According to the *Polaris Project*, each year 21 million people are enslaved worldwide to generate a profit of \$32 billion for their captors. To combat this global problem, organizations like Polaris Project maintain extensive databases of information collected from various public and private sources. The organization reports that it may collect up to 170 different quantitative and qualitative variables per case record, including first-hand data obtained through its National Human Trafficking Resource Center Hotline. In 2013, the NHTRC received 31,945 phone calls, 1,488 e-mails, 1,669 tips from online form submissions, and 787 SMS threads. In order to leverage this vast amount of data, the organization uses the Palantir Gotham analytics platform to track trafficking rings, quickly identify discrete human-trafficking events, and mobilize appropriate response units.



Source: [\*Polaris Project, Launching the Global Human Trafficking Hotline Network\*](#)







