

Identifying Behavioral Health Conditions in Police Encounters

Stanford Data Science for Social Good
August 2022

Our Team



Alex Lerner
DSSG Fellow



Amanda Xu
DSSG Fellow



Grishma Bhattacharai
DSSG Fellow



Salil Goyal
DSSG Fellow



Qian Zhao
Technical Mentor



Dr. Balasubramanian Narasimhan
Faculty Advisor



Lisa Pickoff-White
Community Partner

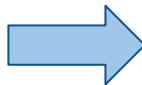
Agenda

- **Project Background**
- **Data Hierarchy**
- **Our Pipeline**
 - Page Type Classification
 - Forms Pipeline
 - Narratives Pipeline
 - DistilBERT Model Justification
 - Text Cleaning
 - Training Labels
 - Model Performance
- **Conclusion and Future Work**

Reporting shows need for law enforcement accountability

California Senate Bill 1421

*made police records relating to officer **use-of-force incidents**, sexual assault and acts of dishonesty **publicly accessible***



NEWS

Bakersfield Police Broke 31 People's Bones in Four Years. No Officer Has Been Disciplined for It

By Lisa Pickoff-White, KQED; Ross Ewald and Danielle Echeverria, Stanford University; Anne Daugherty, UC Berkeley Jun 16, 2021

On Nov. 24, 2017, Robert Cruz Jr. biked north along Baker Street, on a quiet block straddling Bakersfield's once-thriving old town and struggling new, restaurants interspersed with a rehab center and a prepaid phone store. A little before midnight, two officers noticed that the 37-year-old Cruz didn't have a front light on his bicycle. A patrol officer chased Cruz to a nearby yard. There, Cruz crouched behind a child's play tunnel, and the officer struck his arm with a baton.

In California: **3,600+ people seriously injured or killed** by law enforcement officers (2016-2020)

Source: <https://www.kqed.org/news/11878013/bakersfield-police-broke-31-peoples-bones-in-four-years-no-officer-has-been-disciplined-for-it>

Source: <https://www.kvpr.org/local-news/2022-04-12/bakersfield-police-department-fails-to-identify-people-in-crisis-thwarting-reform>

Behavioral health factors require police de-escalation

“Behavioral health conditions”:

- Entails both mental health and substance abuse

Why does it matter?

- Officers are supposed to attempt to de-escalate situations with mentally ill or disabled or intoxicated people, who may not respond to commands in the same way

Bakersfield Police Department fails to identify people in crisis, thwarting reform

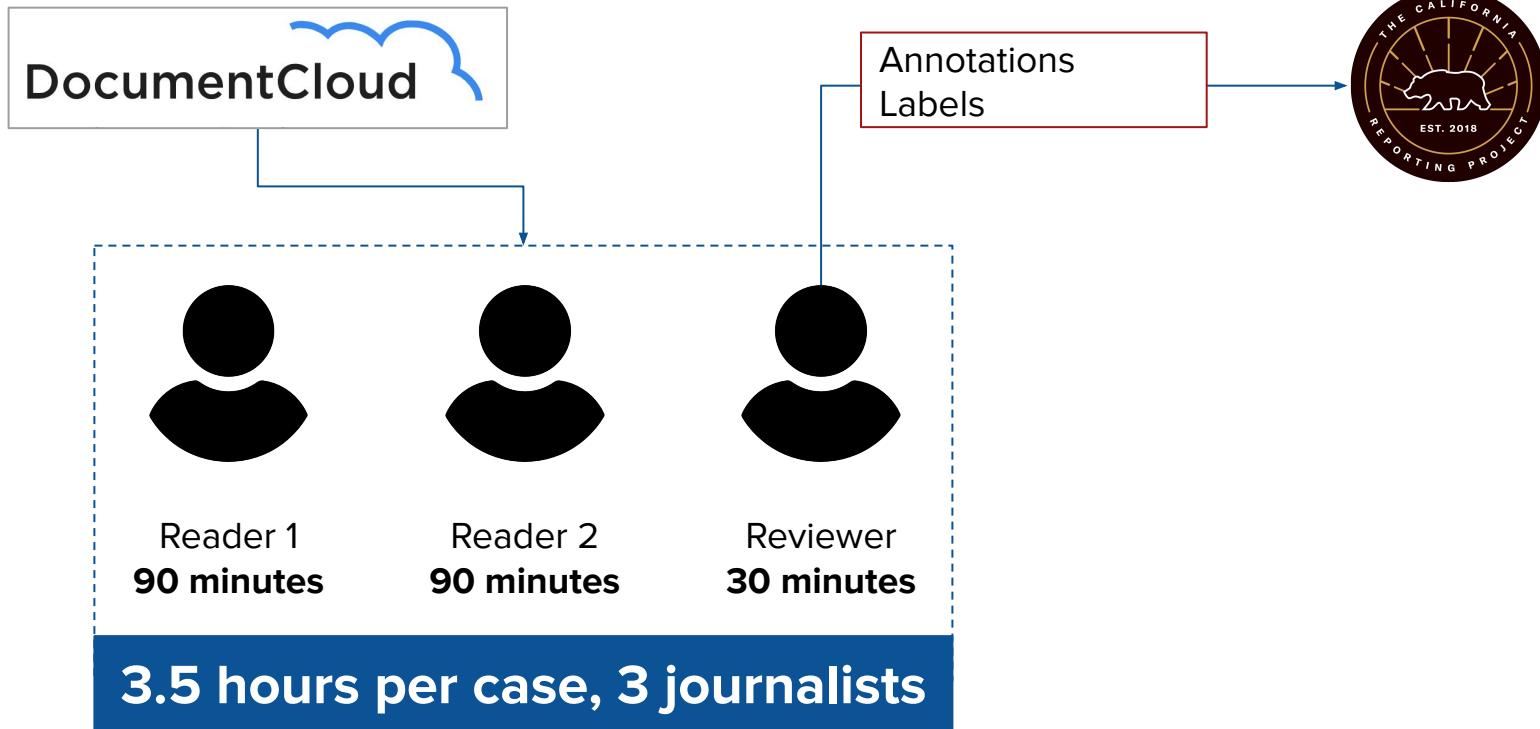
Cases of serious injury/death that involved behavioral health conditions in Bakersfield:

3%
reported rate

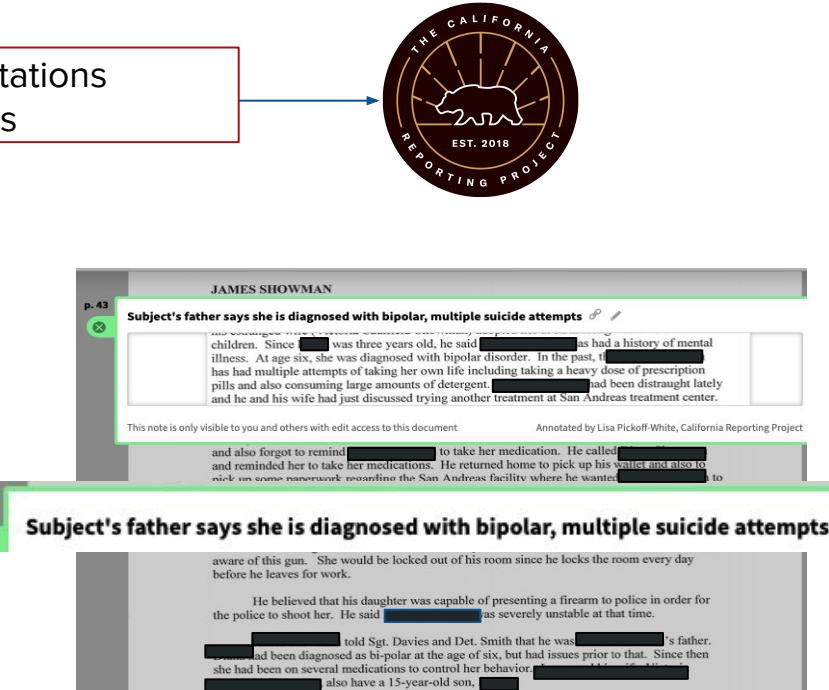
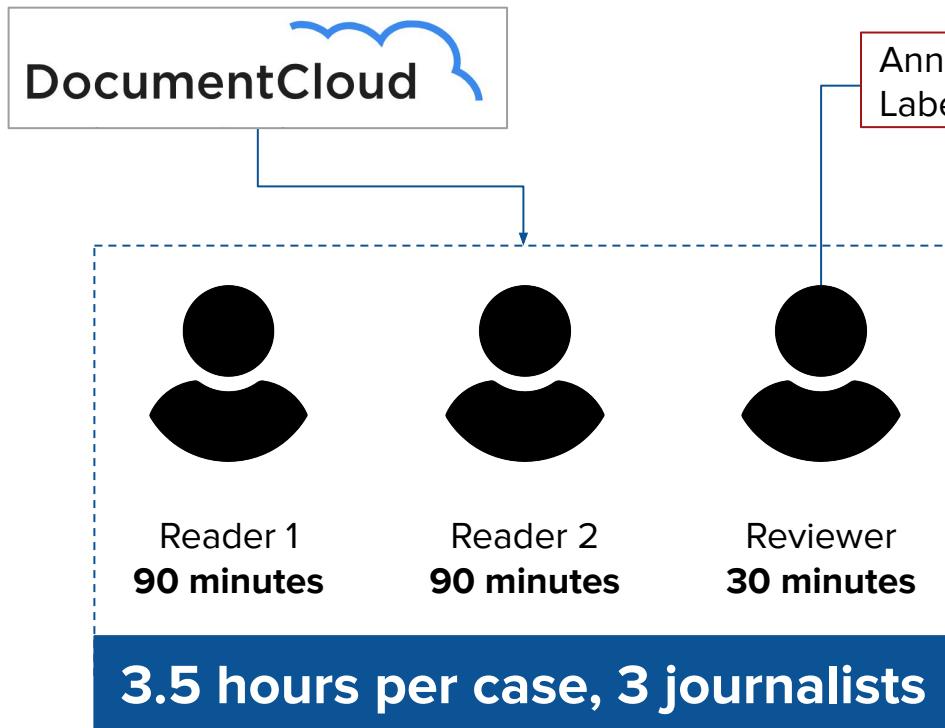
vs.

44%
actual rate

Laborious case review process requires 3.5 man hours



Laborious case review process requires 3.5 man hours



Scale of data creates need for automation

1200 hrs

(That's 30 work weeks!)

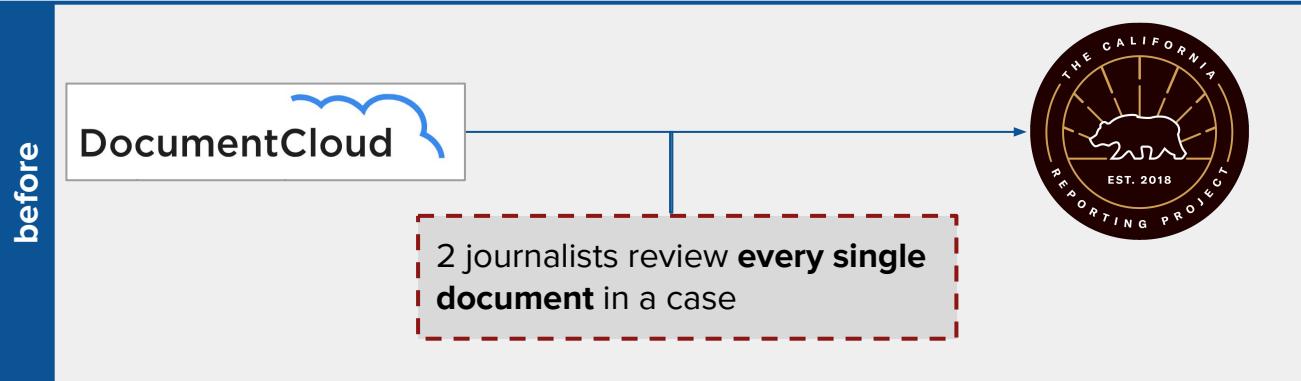
That's how much time*
journalists have spent
on 6 agencies.

There are up to

464

agencies

that are unreviewed.



*Note: 3.5 hours per case (based on BigLocal estimates) * 350 annotated cases = 1,225 hours

Our aim is to reduce workload via automated classification

Project Goal: Automatically classify whether each document contains mention of behavioral health conditions.

unsure if either of them had been hit by the first shot, was worried about being shot at again, and in order to protect his life and his partner's life, returned fire. [REDACTED] was struck twice by the return fire and died.

The Coroner determined [REDACTED] received two gunshot wounds—one to his right arm and another to the lower abdomen which was fatal. A toxicology test of [REDACTED] revealed his blood alcohol content to be a .07%. It also showed that [REDACTED] had cocaine in his system. According to the report, cocaine is a central nervous stimulant drug, and effects from its usage include euphoria, excitement, restlessness, risk taking, and aggression.

Suspect Injuries: (Select only one)

- Not Injured
 - Injured, No Treatment Needed
 - Injured, Refused Treatment
 - Injured, Treated
- Hospital Administrative Clearance ONLY (TASER/Carotid/Projectile Impact)
 - Fatal
 - Unknown, suspect escaped

Signs of Chemical Influence (Drugs and/or Alcohol):

- Yes
- No

Signs of Mental Illness:

- Yes
- No

Was the suspect: (Select all that apply)

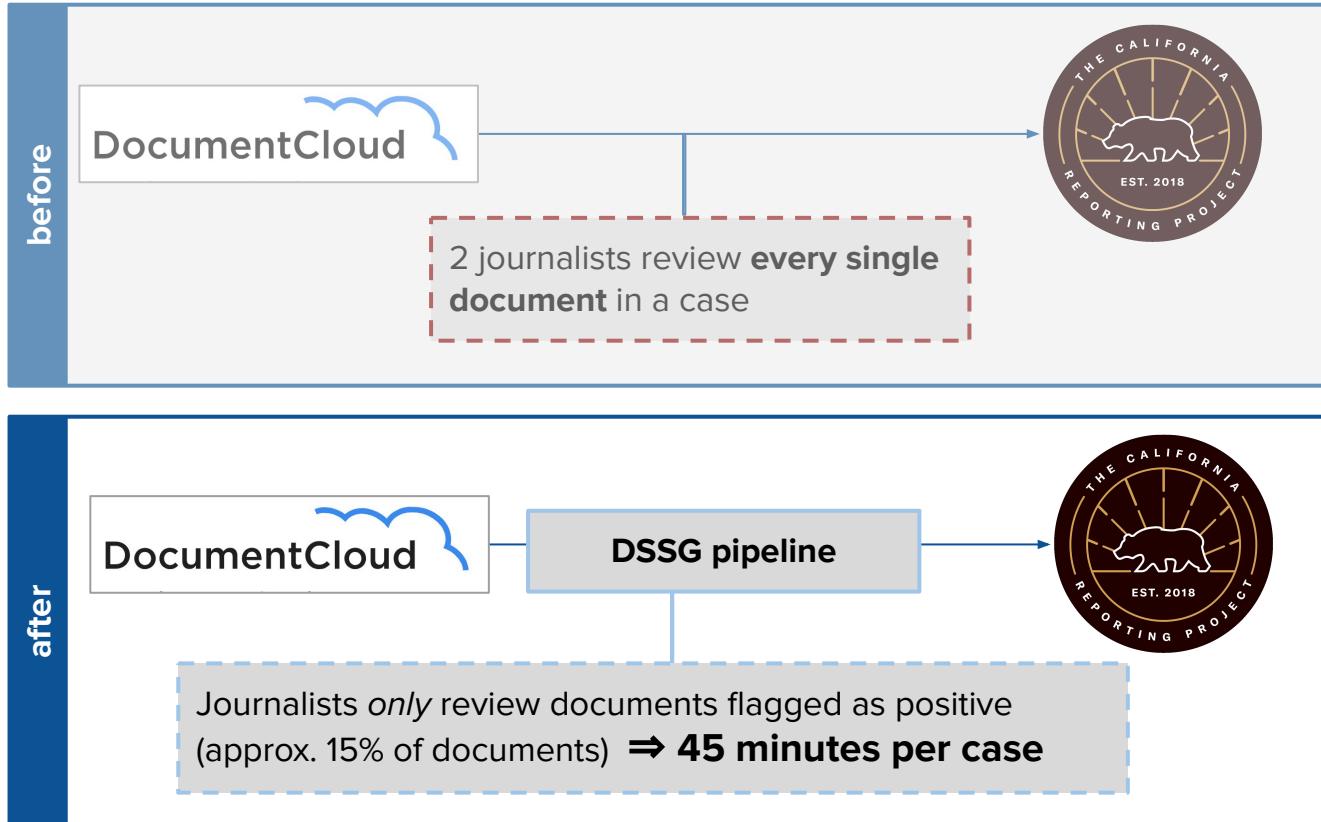
- A fugitive (Warrant, Felony Affidavit, Fleeing Arrest)
- On parole
- On probation

Mentions that subject's BAC was 0.07% and had cocaine in his system
⇒ classify as **True**

Mentions signs of chemical influence (Drugs/Alcohol)
⇒ classify as **True**

Automation reduces man hours required by 75%

75%
less time required
for journalists to
review a case with
automation

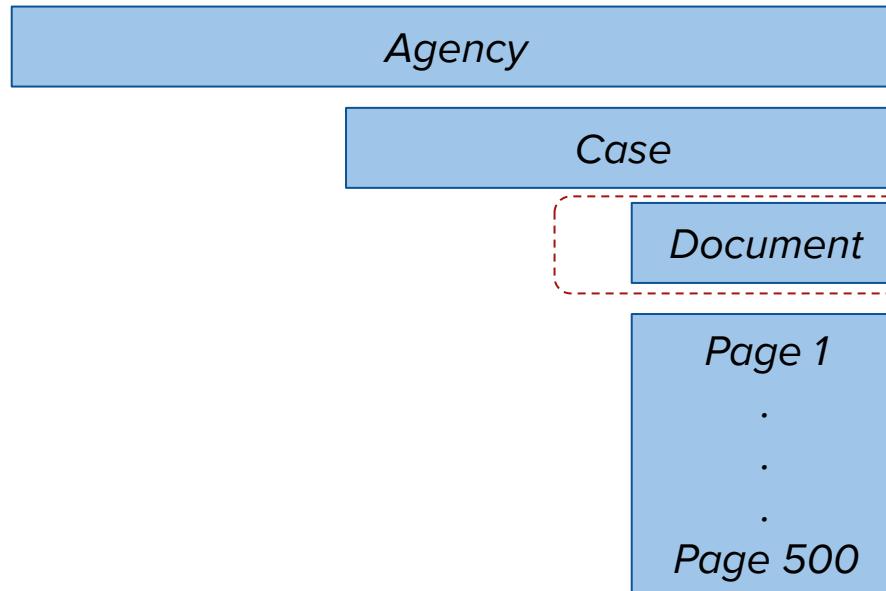


Agenda

- **Project Background**
- **Data Hierarchy**
- **Our Pipeline**
 - Page Type Classification
 - Forms Pipeline
 - Narratives Pipeline
 - DistilBERT Model Justification
 - Text Cleaning
 - Training Labels
 - Model Performance
- **Conclusion and Future Work**

Data hierarchy spans 4 levels, from agency to page

Data hierarchy



Our training data

3 agencies

350+ cases

1600+ documents

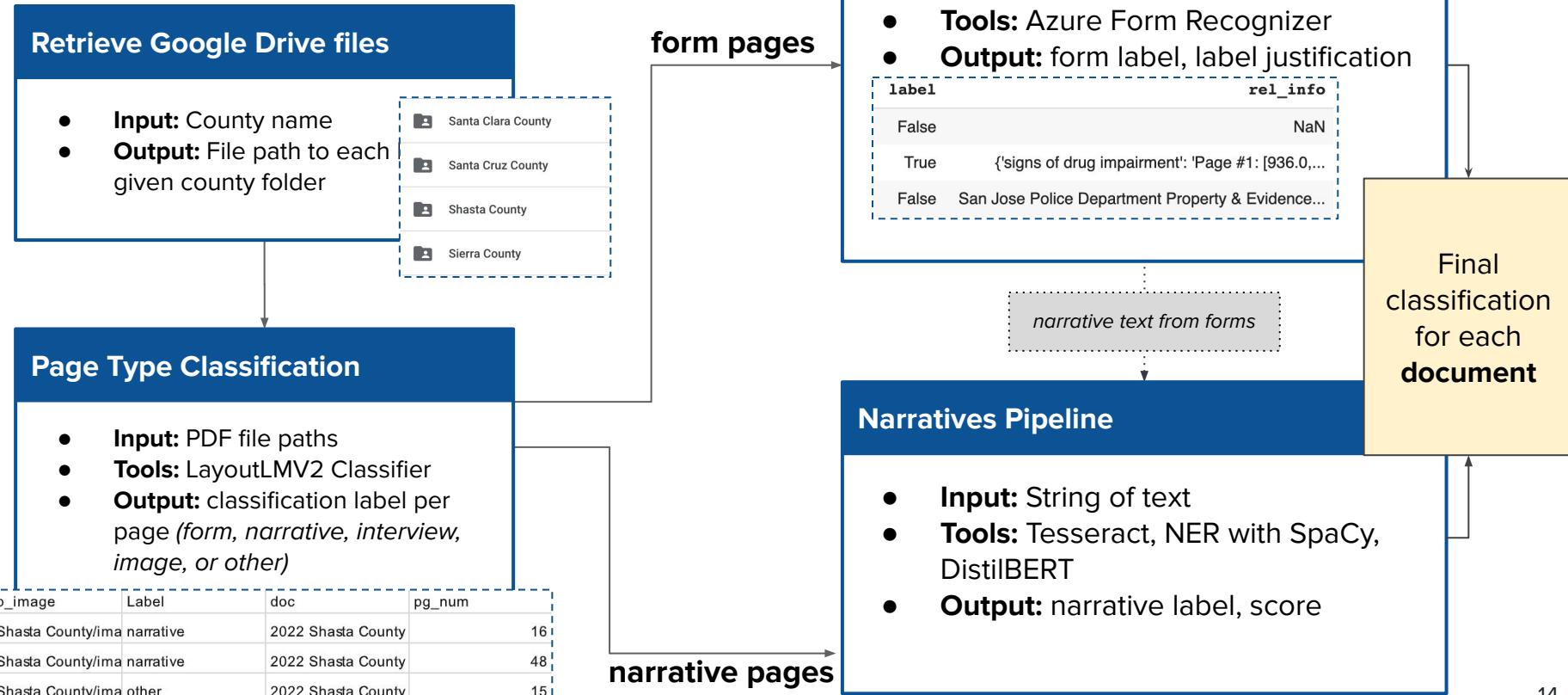
*Each document between
one and several hundred
pages*

Recall:
We want to
classify at the
Document Level!

Agenda

- **Project Background**
- **Data Hierarchy**
- **Our Pipeline**
 - Page Type Classification
 - Forms Pipeline
 - Narratives Pipeline
 - DistilBERT Model Justification
 - Text Cleaning
 - Training Labels
 - Model Performance
- **Conclusion and Future Work**

Pipeline Overview



Page Type Classification

Retrieve Google Drive files

- **Input:** County name
- **Output:** File path to each given county folder

Santa Clara County
Santa Cruz County
Shasta County
Sierra County

form pages

Page Classification

- **Input:** PDF file paths
- **Tools:** LayoutLMV2 Classifier
- **Output:** classification label per page (*form, narrative, interview, image, or other*)

rep_image	Label	doc	pg_num
./Shasta County/ima narrative	narrative	2022 Shasta County	16
./Shasta County/ima narrative	narrative	2022 Shasta County	48
./Shasta County/ima other	other	2022 Shasta County	15

narrative pages

Forms Pipeline

- **Input:** file path to form pages
- **Tools:** Azure Form Recognizer
- **Output:** form label, label justification

label	rel_info
False	NaN
True	{'signs of drug impairment': 'Page #1: [936.0,...
False	San Jose Police Department Property & Evidence...

narrative text from forms

Narratives Pipeline

- **Input:** String of text
- **Tools:** Tesseract, NER with SpaCy, DistilBERT
- **Output:** narrative label, score

Final classification for each document

form

Related Image - FORCE RESPONSE REPORT/USE OF FORCE

Incident Description: Reference Number:

SJPD FORCE RESPONSE REPORT 100-000100000-000-000000000000000000 **IN-534-064R**

Date 08/03/18 Time 14:28 Enforcement Zone 270-Area No.

Suspect Information:

Suspect Name (last, first, middle) DE LA CRUZ, RENE
□ Unknown; suspect named (enter specific name) in name field, i.e., "suspect #1"
DOB 07/14/77 or Approximate Age 41 DOB unknown
Male Female

City of Birth Los Angeles, CA
State of Birth California
Place of Arrest [REDACTED]

Last Known Address Where Suspect Was Located at the Time of the Force Response [REDACTED]

Address 14403 S TELEGRAPH RD
City LOS ANGELES
State CA

Respect Status Select only one
Hostile Hostile Administrative Clearance ONLY
Hostile Reduced Threatened Total
Hostile Unarmed Suspect Unarmed

Signs of Chemical Influence (Drowsy, Dizzy, Headache)
Yes No

Reported Mental Health
Yes No

Are the suspect armed? Yes No
If yes, describe the weapon(s) in the Placing Arrest? Person

Booking charges: Select all the apply that were communicated against an officer only
§ 87(2)(b) § 87(2)(b) § 87(2)(b) § 87(2)(b)

Booking Information:
X Received call for service Left Called in initial report
Initial Report Suspicious Disturbance
Armed Unarmed Armed & Unarmed
Armed on Citizen Armed on Officer Armed on Vehicle
Armed on Bike Armed on Boat Armed on Other
Crime in Progress Only violent

Primary Officer Activity Investigate Prior to Force Response
Investigate Arrest Search Transporting in Vehicle Booking Guide
Arrested Handcuffed Search Handcuffed
Handcuffed Search Handcuffed Search
Gloves worn Search

INCIDENT REPORT PC-0000000000000000000000000000000000

NEVADA COUNTY SHERIFF'S OFFICE 100-0000000000000000000000000000000000

SEARCHED INDEXED SERIALIZED FILED RELEASED

REPORT DATE 08/03/2018 REPORT TIME 14:40:40 BY PC-AUTOMATIC

SEARCHED INDEXED SERIALIZED FILED RELEASED
BY PC-AUTOMATIC ON 08/03/2018 AT 14:40:40

LEATHERMAN STYLE KNIFE
DESCRIPTION Black LOCATION Right Hand FOUND RECOVERED RELEASER

SEARCHED INDEXED SERIALIZED FILED RELEASED
BY PC-AUTOMATIC ON 08/03/2018 AT 14:40:40

E-329 GLOVES
DESCRIPTION Black LOCATION Left Hand FOUND RECOVERED RELEASER

SEARCHED INDEXED SERIALIZED FILED RELEASED
BY PC-AUTOMATIC ON 08/03/2018 AT 14:40:40

E-789 HAMMOCK
DESCRIPTION Black LOCATION Right Side Bed FOUND RECOVERED RELEASER

SEARCHED INDEXED SERIALIZED FILED RELEASED
BY PC-AUTOMATIC ON 08/03/2018 AT 14:40:40

ADDITIONAL INFORMATION

SEARCHED INDEXED SERIALIZED FILED RELEASED
BY PC-AUTOMATIC ON 08/03/2018 AT 14:40:40

100-0000000000000000000000000000000000

SEARCHED INDEXED SERIALIZED FILED RELEASED
BY PC-AUTOMATIC ON 08/03/2018 AT 14:40:40

E-012 UNARMED ELECTRONIC EQUIPMENT
DESCRIPTION Black LOCATION Right Side Bed FOUND RECOVERED RELEASER

SEARCHED INDEXED SERIALIZED FILED RELEASED
BY PC-AUTOMATIC ON 08/03/2018 AT 14:40:40

LOCATED INSIDE A BLACK BAG WITH OTHER CLOTHES AND ELECTRONIC EQUIPMENT

SEARCHED INDEXED SERIALIZED FILED RELEASED
BY PC-AUTOMATIC ON 08/03/2018 AT 14:40:40

E-158 KNIFE
DESCRIPTION Black LOCATION Right Hand FOUND RECOVERED RELEASER

SEARCHED INDEXED SERIALIZED FILED RELEASED
BY PC-AUTOMATIC ON 08/03/2018 AT 14:40:40

E-016 LIGHT STAND
DESCRIPTION Black LOCATION Left Hand FOUND RECOVERED RELEASER

SEARCHED INDEXED SERIALIZED FILED RELEASED
BY PC-AUTOMATIC ON 08/03/2018 AT 14:40:40

ADDITIONAL INFORMATION

SEARCHED INDEXED SERIALIZED FILED RELEASED
BY PC-AUTOMATIC ON 08/03/2018 AT 14:40:40

100-0000000000000000000000000000000000

SEARCHED INDEXED SERIALIZED FILED RELEASED
BY PC-AUTOMATIC ON 08/03/2018 AT 14:40:40

E-003 RETAINER
DESCRIPTION Black LOCATION Right Hand FOUND RECOVERED RELEASER

SEARCHED INDEXED SERIALIZED FILED RELEASED
BY PC-AUTOMATIC ON 08/03/2018 AT 14:40:40

ADDITIONAL INFORMATION

SEARCHED INDEXED SERIALIZED FILED RELEASED
BY PC-AUTOMATIC ON 08/03/2018 AT 14:40:40

narrative

Pursuit/OIS 2/22/18 - Jeff Tyner Page 1 of 1

Pursuit/OIS 2/22/18

Jeff Tyner
Tue 3/02/2018 1609 AM [REDACTED]

Greetings,
I am conducting an Internal Affairs investigation on behalf of the Nevada County Sheriff's Office in reference to the Pursuit which ended in a multiple Officer involved Shooting dated February 22, 2018. My records indicate that you responded to the scene.
I am interested in setting up an interview with you as a witness to the investigation and anticipate that I will be able to contact me at the below number.
Sergeant Jeff Tyner
Nevada County Sheriff's Office
950 Main Street
Nevada City, CA 95959
(530) 265-4500 (Main)
(530) 265-4528 (SRO)
(530) 265-4543 Fax

[REDACTED]

<https://websmail.nevadacounty.ca.us/> 3/20/2018

form

INCIDENT REPORT PC-0000000000000000000000000000000000

NEVADA COUNTY SHERIFF'S OFFICE 100-0000000000000000000000000000000000

SEARCHED INDEXED SERIALIZED FILED RELEASED
BY PC-AUTOMATIC ON 08/03/2018 AT 14:40:40

BROWNSVILLE 100-0000000000000000000000000000000000

SEARCHED INDEXED SERIALIZED FILED RELEASED
BY PC-AUTOMATIC ON 08/03/2018 AT 14:40:40

LEATHERMAN STYLE KNIFE
DESCRIPTION Black LOCATION Right Hand FOUND RECOVERED RELEASER
BY PC-AUTOMATIC ON 08/03/2018 AT 14:40:40

E-329 GLOVES
DESCRIPTION Black LOCATION Left Hand FOUND RECOVERED RELEASER
BY PC-AUTOMATIC ON 08/03/2018 AT 14:40:40

E-789 HAMMOCK
DESCRIPTION Black LOCATION Right Side Bed FOUND RECOVERED RELEASER
BY PC-AUTOMATIC ON 08/03/2018 AT 14:40:40

ADDITIONAL INFORMATION

SEARCHED INDEXED SERIALIZED FILED RELEASED
BY PC-AUTOMATIC ON 08/03/2018 AT 14:40:40

100-0000000000000000000000000000000000

SEARCHED INDEXED SERIALIZED FILED RELEASED
BY PC-AUTOMATIC ON 08/03/2018 AT 14:40:40

E-012 UNARMED ELECTRONIC EQUIPMENT
DESCRIPTION Black LOCATION Right Side Bed FOUND RECOVERED RELEASER
BY PC-AUTOMATIC ON 08/03/2018 AT 14:40:40

LOCATED INSIDE A BLACK BAG WITH OTHER CLOTHES AND ELECTRONIC EQUIPMENT

SEARCHED INDEXED SERIALIZED FILED RELEASED
BY PC-AUTOMATIC ON 08/03/2018 AT 14:40:40

E-158 KNIFE
DESCRIPTION Black LOCATION Right Hand FOUND RECOVERED RELEASER
BY PC-AUTOMATIC ON 08/03/2018 AT 14:40:40

E-016 LIGHT STAND
DESCRIPTION Black LOCATION Left Hand FOUND RECOVERED RELEASER
BY PC-AUTOMATIC ON 08/03/2018 AT 14:40:40

ADDITIONAL INFORMATION

SEARCHED INDEXED SERIALIZED FILED RELEASED
BY PC-AUTOMATIC ON 08/03/2018 AT 14:40:40

100-0000000000000000000000000000000000

SEARCHED INDEXED SERIALIZED FILED RELEASED
BY PC-AUTOMATIC ON 08/03/2018 AT 14:40:40

E-003 RETAINER
DESCRIPTION Black LOCATION Right Hand FOUND RECOVERED RELEASER
BY PC-AUTOMATIC ON 08/03/2018 AT 14:40:40

ADDITIONAL INFORMATION

SEARCHED INDEXED SERIALIZED FILED RELEASED
BY PC-AUTOMATIC ON 08/03/2018 AT 14:40:40

Page Type Classification

Goal: Separate forms from narratives pages within a document

Input

file path to Google Drive

Hugging Face

Page type classifier*

Output

document, page number, page path

page label

*created by Hellina Hailu Nigatu at UC Berkeley

Forms Pipeline

Retrieve Google Drive files

- **Input:** County name
- **Output:** File path to each given county folder

Santa Clara County
Santa Cruz County
Shasta County
Sierra County

form pages

Page Classification

- **Input:** PDF file paths
- **Tools:** LayoutLMV2 Classifier
- **Output:** classification label per page (*form, narrative, interview, image, or other*)

rep_image	Label	doc	pg_num
./Shasta County/ima narrative		2022 Shasta County	16
./Shasta County/ima narrative		2022 Shasta County	48
./Shasta County/ima other		2022 Shasta County	15

narrative pages

Forms Pipeline

- **Input:** file path to form pages
- **Tools:** Azure Form Recognizer
- **Output:** form label, label justification

label	rel_info
False	NaN
True	{'signs of drug impairment': 'Page #1: [936.0,...
False	San Jose Police Department Property & Evidence...

Final classification for each document

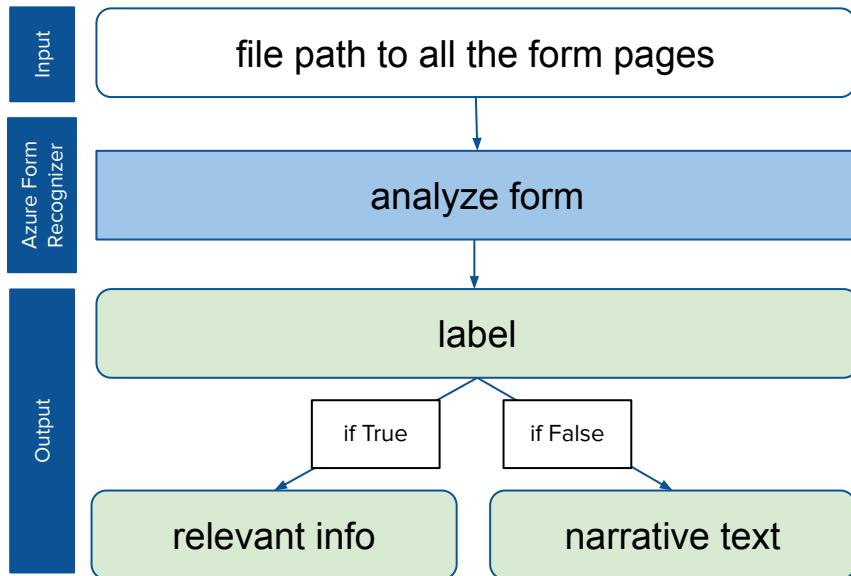
narrative text from forms

Narratives Pipeline

- **Input:** String of text
- **Tools:** Tesseract, NER with SpaCy, DistilBERT
- **Output:** narrative label, score

Forms Pipeline

Goal: Find forms that mention behavioral health conditions + explain why



For each form page...

1. Analyze checkboxes. Select the ones that are:
 - Selected
 - Relevant (contains some relevant word; e.g., “intoxication”, “influence”)

Return True, checkbox title/location

2. If no relevant checkboxes → retrieve any narrative text on page

Return False, narrative text

Example outputs of forms

Example of Positive Label

How many suspects were present when the officer responded with force?

Single Suspect Multiple Suspects

Suspect Injuries: (Select only one)

Not Injured Hospital Administrative Clearance ONLY (TASER/Carotid/Projectile Impact)
 Injured, No Treatment Needed Fatal
 Injured, Refused Treatment Unknown, suspect escaped
 Injured, Treated

Signs of Chemical Influence (Drugs and/or Alcohol):

Yes No

Signs of Mental Illness:

Yes No

Was the suspect: (Select all that apply)

A fugitive (Warrant, Felony Affidavit, Fleeing Arrest) On parole On probation



True

```
{'Signs of Chemical Influence (Drugs and/or Alcohol)'::  
'Page #1: [313.0, 904.0], [886.0, 904.0], [886.0, 930.  
0], [313.0, 930.0]'}
```

Example of Negative Label

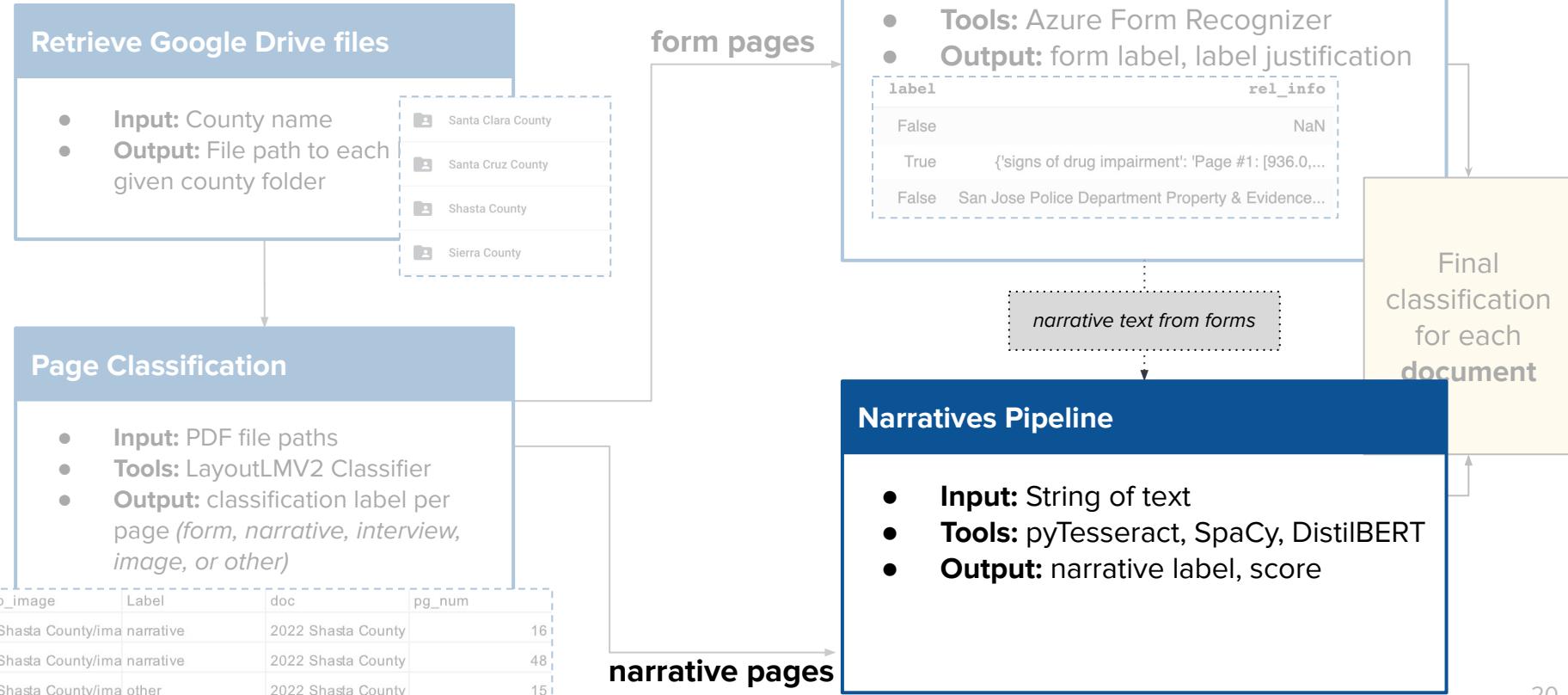
Rep Unit	0266	District	LA	Sector	LW	From Date	02/02/2019	From Time	15:26
Members	Assignment								
Entered By	2860	Assignment	FIELD SERVICE SEAPORT 4TH WATCH PHASE 2	Role	Not Transfer	Prop Trans Stat	Successful	Prop Trans Stat	Successful
Approving Officer	2052	Approval Date	02/07/2019	Approval Time	13:09:02				
Summary Narrative Supplemental report - Scene perimeter BWC activated Modus Operandi Crime Code(s) Assaults Narrative									
NOTIFICATION: On 02-02-19 at 1537 hours, Officer Hammer (B7) and I (Gushchin/B8), were dispatched to assist officers on a scene of an officer involved shooting located at [REDACTED]									
INVESTIGATION: On arrival to the scene I signed the Crime Scene Log. I was advised by officers on scene to circulate the North side of the perimeter for possible witnesses. As I was checking the area, I noticed a black male adult approach officers on scene of the East side of the perimeter. The black male adult was bald and he was not wearing a shirt. The black male appeared agitated due to him hitting his arms together and due to him constantly shouting, "That's my little sister in there," "Where's that nigga at," "I'm going to fuck him up." The black male eventually entered a newer model, dark gray or black 4 door Hyundai Elantra or Hyundai Sonata, and left the area.									



False,

```
'On arrival to the scene I signed the Crime Scene Log. I was advised by officers on scene to circulate the North side of the perimeter for possible witnesses. As I was checking the area, I noticed a black male adult approach officers on scene of the East side of the perimeter. The black male adult was bald and he was not wearing a shirt. The black male appeared agitated due to him hitting his arms together and due
```

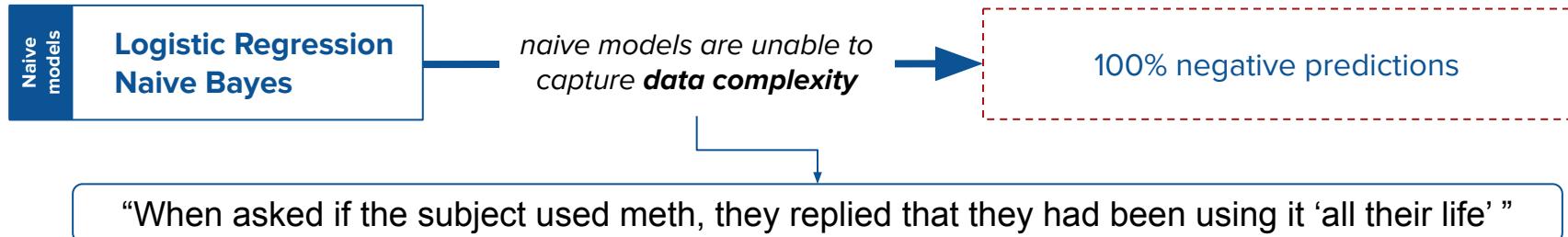
Narratives Pipeline



Agenda

- **Project Background**
- **Data Hierarchy**
- **Our Pipeline**
 - Page Type Classification
 - Forms Pipeline
 - Narratives Pipeline
 - DistilBERT Model Justification
 - Text Cleaning
 - Training Labels
 - Model Performance
- **Conclusion and Future Work**

Complex model needed to learn language meaning



Recent advances in Natural Language Processing (NLP) have championed the use of Deep Neural Models which have shown tremendous performance on a variety of tasks

- Question Answering, Summarization, Translation, and also Text Classification

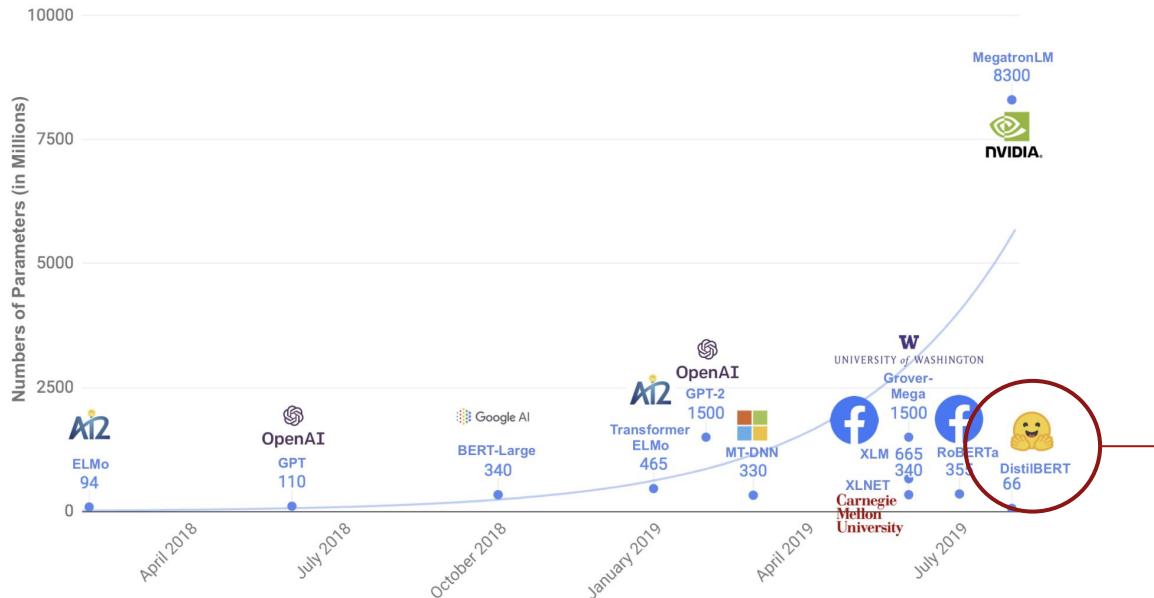
Training our model



Pre-trained Transformer model fits our needs

- In 2017, Google released a new model architecture that became the model of choice for many NLP tasks: the Transformer
- Processes sequences at a time
- Utilizes a “self-attention” mechanism that can weight different parts of each sequence
- Is pre-trained on billions of texts from the web and Wikipedia, and can be fine-tuned for a variety of downstream tasks

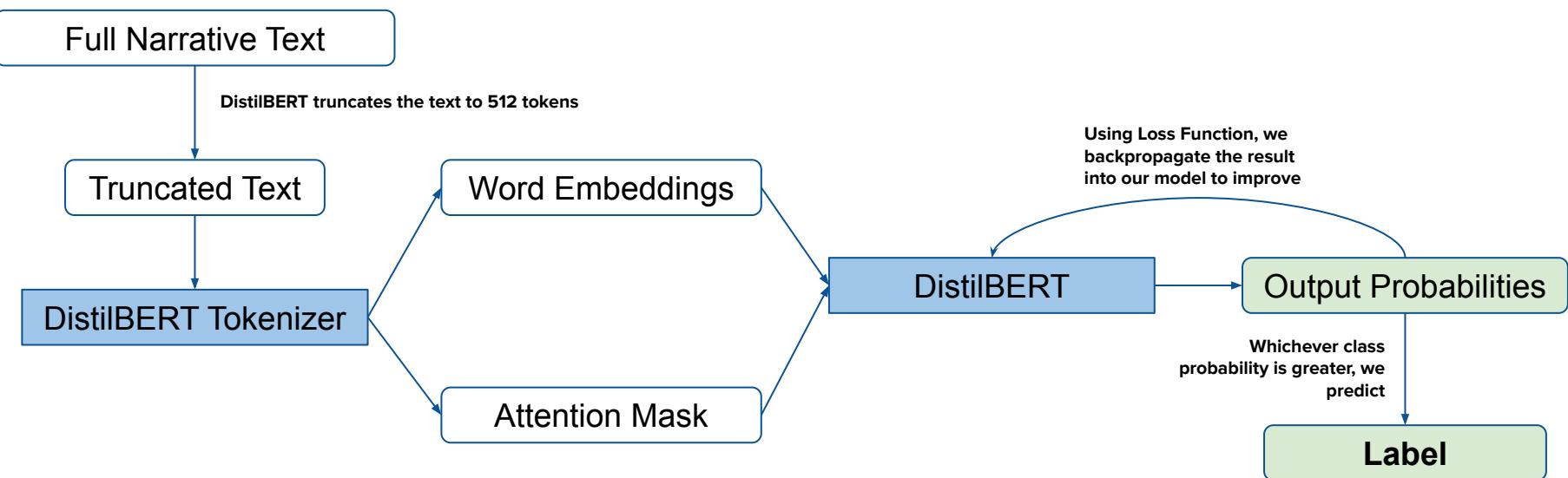
DistilBERT balances training simplicity and performance



Our choice model, DistilBERT, is a variant of Google's Transformer model—but mimics performance with far fewer parameters

⇒ faster training + faster inference, matching our need to classify documents at scale

DistilBERT Pipeline



Agenda

- **Project Background**
- **Data Hierarchy**
- **Our Pipeline**
 - Page Type Classification
 - Forms Pipeline
 - Narratives Pipeline
 - DistilBERT Model Justification
 - Text Cleaning
 - Training Labels
 - Model Performance
- **Conclusion and Future Work**

pyTesseract used to extract text from narrative pages

page	label	link
1	narrative	https://www.documentcloud.org/documents/21067706-15-315-0038-ia-report_redacted
17	image	https://www.documentcloud.org/documents/21067706-15-315-0038-ia-report_redacted
6	narrative	https://www.documentcloud.org/documents/21067706-15-315-0038-ia-report_redacted
12	narrative	https://www.documentcloud.org/documents/21067706-15-315-0038-ia-report_redacted
4	narrative	https://www.documentcloud.org/documents/21067706-15-315-0038-ia-report_redacted
5	narrative	https://www.documentcloud.org/documents/21067706-15-315-0038-ia-report_redacted
15	narrative	https://www.documentcloud.org/documents/21067706-15-315-0038-ia-report_redacted
23	form	https://www.documentcloud.org/documents/21067706-15-315-0038-ia-report_redacted

Output of page type classifier

Narrative page

Statement of Suspect [REDACTED] Vong:

At the very beginning of the interview, [REDACTED] spontaneously (without being asked) said, "I'm just glad I didn't kill anybody."

Wong stated he lives by himself in the townhouse located at [REDACTED] Bautista Place. He gave the police consent to search his residence and the vehicle in the garage. His [REDACTED], whom he identified as [REDACTED] [REDACTED] had lived in the residence with him until they broke up 3-4 months ago. Wong does not know where [REDACTED] is currently living and has not seen her for a few months. Wong has a lot of guns because he likes to collect them. He has never been in the military. He works as a nurse at the [REDACTED]. He has two assault rifles (which he built), a shotgun, two HK 9 mm handguns, and a Sig Sauer 9 mm handgun. He usually keeps these guns in various safes in his residence.

All of his guns are registered in California under his name. He said that the Sig Sauer 9 mm gun that police officers found in the bushes belonged to him and he dropped it from the roof just prior to the incident with the police. When Sergeant Hamblin asked him what he was doing on the roof he said, "I don't know." He said he had one of the assault rifles slung across his chest and a 1911 handgun strapped to his chest when he was contacted by the police at the front of his residence.

pyTesseract

Input: Image of page

Output: Text extracted
from the image

Raw text after extraction (pyTesseract output)

tement. Refer to the General Offense report for additional information. Statement of Suspect [REDACTED]: At the very beginning of the interview, [REDACTED] spontaneously (without being asked) said, "I'm just glad I didn't kill anybody." [REDACTED] stated he lives by himself in the townhouse located at Bautista Place. He gave the police consent to search his residence and the vehicle in the garage. His [REDACTED], whom he identified as [REDACTED] [REDACTED] had lived in the residence with him until they broke up 3-4 months ago. Wong does not know where [REDACTED] is currently living and has not seen her for a few months. [REDACTED] has a lot of guns because he likes to collect them. He has never been in the military. He works as a nurse at the [REDACTED]. He has two assault rifles (which he built), a shotgun, two HK 9 mm handguns, and a Sig Sauer 9 mm handgun. He usually keeps these guns in various safes in his residence. All of his guns are registered in California under his name. He said that the Sig Sauer 9 mm gun that police officers found in the bushes belonged to him and he dropped it from the roof just prior to the incident with the police. When [REDACTED] asked him what he was doing on the roof he said, "I don't know." He said he had one of the assault rifles slung across his chest and a 1911 handgun strapped to his chest when he was contacted by the police at the front of his residence.

Named Entity Recognition (with spaCy)

NER classifies named entities mentioned in unstructured text into predefined categories (e.g., person names, organizations, locations, medical codes, time expressions, quantities, monetary values, percentages, etc.)

Raw Text

'KERN COUNTY DISTRICT ATTORNEY\n\nO F F I C E\n\nM E M O R A N D U M\n\nTO: District Attorney Lisa Green\nFROM: DDA R
AM\nDATE: January 20, 2017\nSUBJECT: OIS near Feliz Dr; BPD Officers Shaff, Mayberry, Mueller, Pace, and Barber; Sus
pect [REDACTED] [REDACTED] [REDACTED] [REDACTED]\nBrief Factual Synopsis:\nOn August 18, 2016, just before 3:30am, [REDACTED] [REDACTED] [REDACTED] [REDACTED] ra
from a car that had just ended a lengthy vehicle pursuit with BPD. He told the occupants in the car with him that
he was "not going back to prison." [REDACTED] [REDACTED] [REDACTED] [REDACTED] armed with a loaded .380 semi-automatic handgun.\nOfficer Ary saw
Villarreal with a gun, but as he attempted to chase him another male passenger exited the car and assaulted Ary. Ary
yelled "gun" but was not heard by the other officers. Five officers were in pursuit [REDACTED] foot through an
apartment complex that is between Feliz Dr and McNew Ct. None of the officers in pursuit knew Villarreal was armed.
As Officer Pace [REDACTED] [REDACTED] [REDACTED] [REDACTED] pointed his gun behind him and shot at her from roughly 10 feet away.\nShe yelled "gun" but was not heard by the other officers as they were approaching Villarreal\nfrom different directions. Other officers mistook the sound of the gunshot for Officer Pace's\ntaser being deployed. Officer Shaff then closed in and began to attempt to take [REDACTED] to the ground. [REDACTED] [REDACTED] [REDACTED] [REDACTED] pointed the gun at Shaff's head and fired. Shaff dropped to the ground in an effort to avoid being shot. Officers saw Villarreal attempt to shoot Shaff. All five officers then shot at [REDACTED] various locations. [REDACTED] continued to run as they shot at him for 15 to 20 feet and then fell. Officers Shaff, Pace, and Mueller approached Villarreal from behind as he was lying on the ground with his back to them. Officer Barber and Mayberry came from the side to approach [REDACTED]. [REDACTED] Barber came within five to ten feet [REDACTED] raised his arm and shot at Barber's head. Barber dropped to the ground. He and the other four officers [REDACTED] After the shooting stopped, Officer Shaff walked ahead of the other officers and as he did he [REDACTED] his arm and body "twitching" in a similar movement to when he fired at [REDACTED]

List of Named Entities (NE) in the raw text

[REDACTED] PERSON	
DDA RAM ORG	
January 20, 2017 DATE	
OIS ORG	
Feliz Dr PERSON	
BPD Officers Shaff ORG	
Mayberry PERSON	
Mueller PERSON	spacy.explain("ORG")
Barber PERSON	
Suspect F PERSON	'Companies, agencies, institutions, etc.'
August 18, 2016 DATE	
3:30am CARDINAL	
[REDACTED] PERSON	
BPD ORG	
[REDACTED] PERSON	
[REDACTED] PERSON	
Ary PERSON	
Five CARDINAL	
spacy.explain("CARDINAL")	
'Numerals that do not fall under another type'	

6 of 18 entities are useful, remaining 12 are masked

TO: District Attorney [REDACTED] in PERSON
FROM: DDA RAM ORG
DATE: January 20, 2017 DATE
SUBJECT: OIS ORG near [REDACTED] PERSON ; BPD Officers Shaff ORG , [REDACTED] PERSON , [REDACTED] PERSON
[REDACTED] PERSON ; Suspect [REDACTED] ORG [REDACTED]

Brief Factual Synopsis:

Noisy Named Entities

Toxicology Report

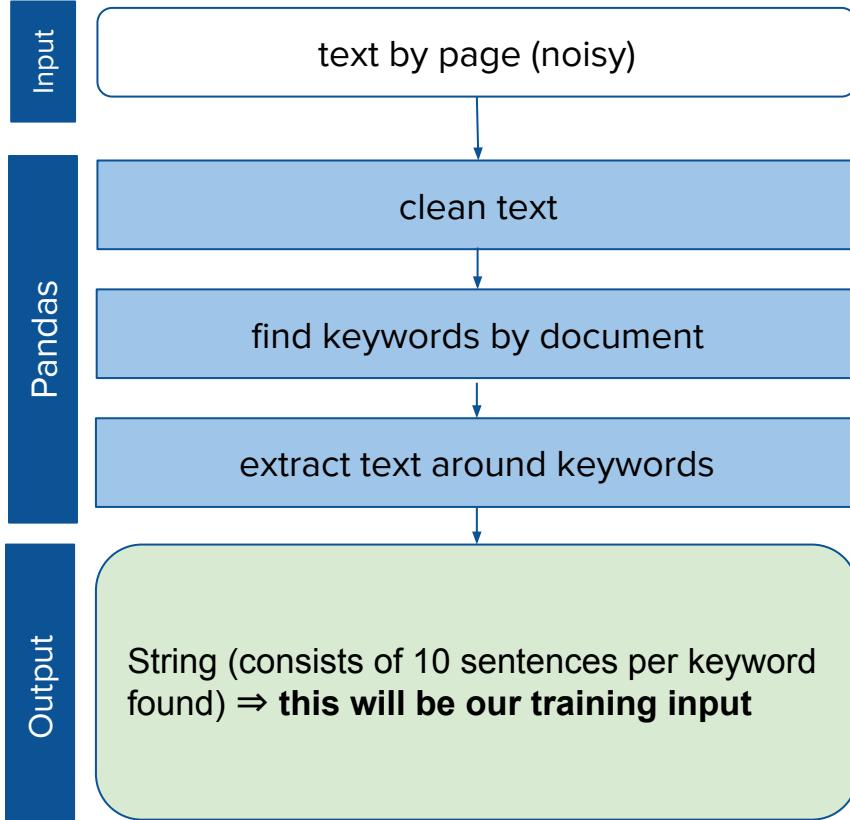
According to the Toxicology Report ORG from NMS Labs ORG , [REDACTED] PERSON tested positive for several substances. He had a small amount of alcohol with a BAC ORG of .019% PERCENT . He tested positive for marijuana. He also tested positive for methamphetamine with 2400 CARDINAL nanograms of methamphetamine in his blood

Relevant Named Entities

After exploring NER, we found that for our purposes:

- 6 categories: ['ORG', 'CARDINAL', 'ORDINAL', 'PRODUCT', 'PERCENT', 'QUANTITY'] may have relevant information
- 12 of 18 categories: noise
- Masked the noisy 12, kept the useful 6 as is

3 text preprocessing steps to prepare text for input to model



NARRATIVE

On 09/09/2019 at approximately 1115 hours, I was on duty with the Richmond Police Department. I was dressed in full Richmond class "D" uniform with a badge on my chest which stated "Police Officer." I was driving vehicle 51-79, a fully marked black and white patrol SUV. I had just finished participating in a pedestrian crosswalk operation to address traffic concerns in the Richmond Hilltop area.

I was listening to the debrief of the operation in the parking lot of the Hilltop Mall near Shane Drive. A US Marshal (Boydland, Lawell) flagged us down in an unmarked vehicle. He stated they were looking for a PAL (Boydland, Lawell) who was driving a Chevrolet Malibu with a license plate CA [REDACTED]. I told him the vehicle was last seen in the Hilltop area near San Pablo Avenue. I ran a check on Boydland and confirmed that he was a wanted felon with a warrant in CLETS. I noticed that the warrant had several "cautions" listed for Boydland to include "violent tendencies" and "knows [REDACTED]". I informed him that [REDACTED] was a [REDACTED] with average build. I saw his CLETS readout had him at height 5'11" and weight 180 lbs.

I began circulating the area and looked through the storefront parking lot near Richmond River and Atlas Drive. I drove onto Quarendon Drive and turned onto Waverly Drive. Approaching San Pablo Ave, my attention was drawn to gray car parked facing eastbound with the stop lamps illuminated. I saw it was only occupied by the driver. As I got closer, I noticed it was a Chevrolet Malibu and the license plate was a match.

Keyword: “Drugs”

Output: 5 sentences

before keyword appears, 5 sentences after keyword appears

Reduced text from **125 sentences → 10 sentences**

KEYWORDS

```
['erratic', 'alcohol', 'mental', 'influence', 'strange', 'odd', 'drug',  
'narcotics', 'crazy', 'toxicology', 'bizarre', 'PCP', 'agitated',  
'5150', 'drunk', 'cocaine', 'involuntary', 'psychiatric', 'methamphetamine',  
'amphetamine', 'marijuana', 'narcotic', 'impairment', 'impaired', 'substance']
```

Agenda

- **Project Background**
- **Data Hierarchy**
- **Our Pipeline**
 - Page Type Classification
 - Forms Pipeline
 - Narratives Pipeline
 - DistilBERT Model Justification
 - Text Cleaning
 - Training Labels
 - Model Performance
- **Conclusion and Future Work**

Existing training data is noisy due to **case-level labels**

We were given **case-level labels**.

Case ID	document	actual label	label we have
080622	incident_report	False	Case Label = True
080622	tox_report	False	Case Label = True
080622	witness_summary	True	Case Label = True

Case Label: True

Problem: noisy labels

Some documents in a case may not mention behavioral health conditions, but the case will still be labeled True

⇒ Need document-level labels for document-level classifications

Annotations offer a way to create document level labels

We were given case-level labels.

Case ID	document	actual label	label we have
080622	incident_report	False	Case Label = True
080622	tox_report	False	Case Label = True
080622	witness_summary	True	Case Label = True

Case Label: True

Problem: noisy labels

Some documents in a case may not mention behavioral health conditions, but the case will still be labeled True

We decided to generate new document-level labels using annotations for our training data.

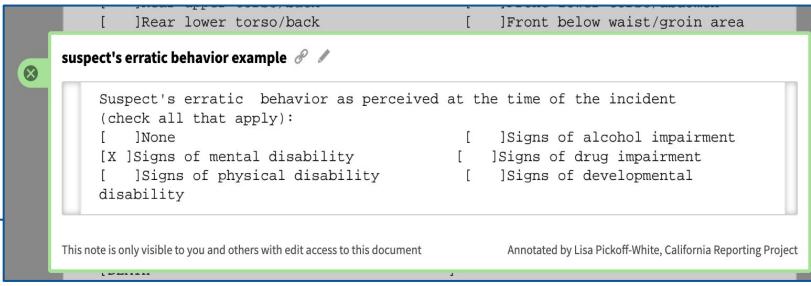
The screenshot shows a portion of a document titled "suspect's erratic behavior example". It contains two checkboxes at the top: "[]Rear lower torso/back" and "[]Front below waist/groin area". Below this is a section titled "suspect's erratic behavior example" with a pencil icon. It contains the text: "Suspect's erratic behavior as perceived at the time of the incident (check all that apply):". There are four groups of checkboxes: 1) "[]None", "[]Signs of alcohol impairment", "[]Signs of drug impairment", "[]Signs of developmental disability". 2) "[X] Signs of mental disability", "[]Signs of physical disability". At the bottom of the form, there is a note: "This note is only visible to you and others with edit access to this document" and "Annotated by Lisa Pickoff-White, California Reporting Project".

Example of an Annotated File on DocumentCloud

Case ID	document	annotation_label
080622	incident_report	True if (relevant annotation), else False
080622	tox_report	True if (relevant annotation), else False

Annotations as labels

We decided to generate new document-level labels using annotations for our training data.

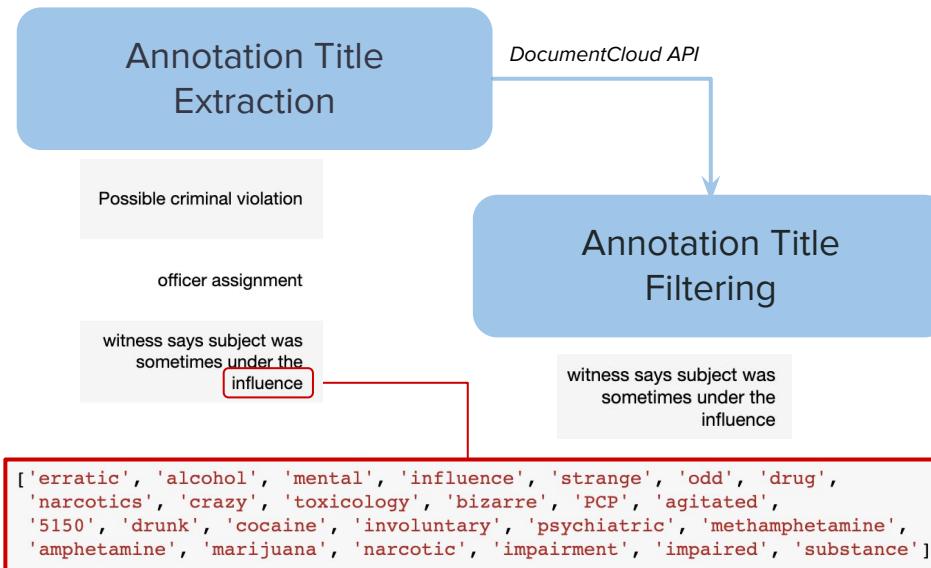


The screenshot shows a document annotation interface. At the top, there are two checkboxes: "[]Rear lower torso/back" and "[]Front below waist/groin area". Below them is a note titled "suspect's erratic behavior example" with a green 'X' icon. The note contains the following text:
Suspect's erratic behavior as perceived at the time of the incident
(check all that apply):
[]None []Signs of alcohol impairment
[x]Signs of mental disability []Signs of drug impairment
[]Signs of physical disability []Signs of developmental
disability

At the bottom of the note, it says "This note is only visible to you and others with edit access to this document" and "Annotated by Lisa Pickoff-White, California Reporting Project". A blue rectangular box points from the left towards the note. To the right of the note, the word "KEYWORDS" is centered above a red-bordered box containing a list of words:

```
[ 'erratic', 'alcohol', 'mental', 'influence', 'strange', 'odd', 'drug',  
'narcotics', 'crazy', 'toxicology', 'bizarre', 'PCP', 'agitated',  
'5150', 'drunk', 'cocaine', 'involuntary', 'psychiatric', 'methamphetamine',  
'amphetamine', 'marijuana', 'narcotic', 'impairment', 'impaired', 'substance' ]
```

Annotation relevance determined based on title content



KEYWORDS

The screenshot shows a document interface with two examples of annotations:

Relevant Annotation: witness says subject was sometimes under the influence

This annotation is highlighted with a green box. The text in the box reads: "he had only seen him under the influence. I asked GARCIA if he was arrested and incarcerated four months ago. I asked GARCIA".

Irrelevant Annotation: Date [REDACTED] penalty imposed

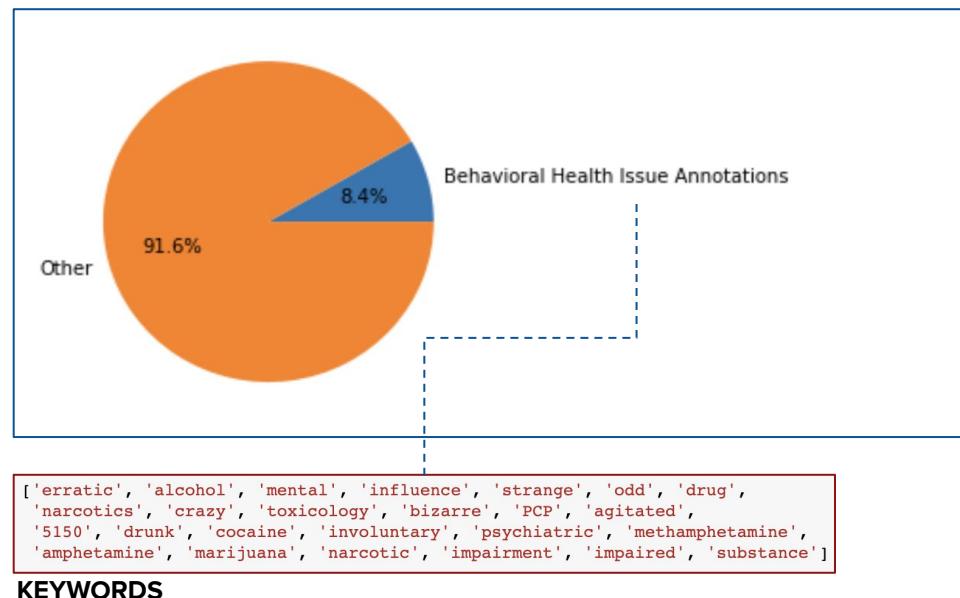
This annotation is highlighted with a red box. The text in the box reads: "the harm caused and/or attempted to conceal the harm through action or inaction." Below this is a signature and the date "5/29".

Irrelevant Annotation: describes administrative detail

Ratio of relevant to irrelevant annotations

8.4%

of annotations proved “relevant,” or had words in their title belonging to our list of keywords



Final train/test set comprises 936 documents

- **Documents:** Originally 766 documents, oversampled positive documents to result in **936** total
- **Text:** Narrative text extracted from document classification, cleaned by NER and then filtered by keyword search
- **Labels:** Presence of useful annotations as a label
- **Test Set:** Randomly sampled 10% of the total dataset to use for evaluation

link	label	text	pages_with_keywords
https://www.documentcloud.org/documents/205175...	0	BAKERSFIELD POLICE DEPARTMENT PRESS RELEASE Ch...	[]
https://www.documentcloud.org/documents/205176...	0	BAKERSFIELD POLICE DEPARTMENT PRESS RELEASE Ch...	[]
https://www.documentcloud.org/documents/205189...	0	They could not find it. said ok you got me and...	[1, 2, 17]
https://www.documentcloud.org/documents/205189...	0	BAKERSFIELD POLICE DEPARTMENT PRESS RELEASE Ch...	[]
https://www.documentcloud.org/documents/205936...	1	immediately drew department issued firearm and...	[10, 20]
https://www.documentcloud.org/documents/206029...	1	See Oficer Beatties supplemental report for fu...	[7, 8, 11]
https://www.documentcloud.org/documents/206045...	1	i stated observed the vehicle flee but was una... [1, 5, 15, 18, 20, 28, 29, 35, 41, 42]	

Agenda

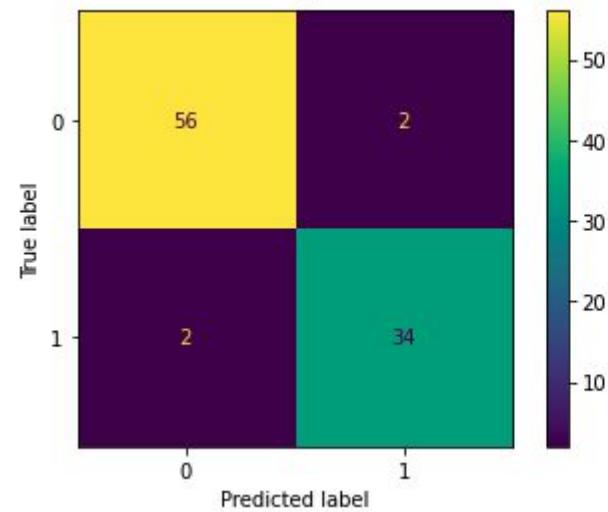
- **Project Background**
- **Data Hierarchy**
- **Our Pipeline**
 - Page Type Classification
 - Forms Pipeline
 - Narratives Pipeline
 - DistilBERT Model Justification
 - Text Cleaning
 - Training Labels
 - Model Performance
- **Conclusion and Future Work**

Narrative model yields strong test accuracy at 96%

Using our document-by-document dataset, we were able to get performance of **~96% accuracy** on our test data (94 documents)!

	precision	recall	f1-score
0	0.97	0.97	0.97
1	0.94	0.94	0.94
accuracy	0.96		

Accuracy: 0.957
Precision: 0.944
Sensitivity/Recall: 0.944
Specificity: 0.965



Final Label Aggregation

Retrieve Google Drive files

- Input: County name
- Output: File path in given county folder

Santa Clara County

form pages

Page Classification

- Input: PDF file path
- Tools: LayoutLMV
- Output: classification label per page (*form, narrative, interview, image, or other*)

rep_image	Label	doc	pg_num
./Shasta County/ima narrative		2022 Shasta County	16
./Shasta County/ima narrative		2022 Shasta County	48
./Shasta County/ima other		2022 Shasta County	15

Forms Pipeline

- Input: file path to form pages
- Tools: Azure Form Recognizer
- Output: form label, label justification

label	rel_info
False	NaN 1: [936.0,... Evidence...]

True if either forms or narratives pipeline output True

Along with relevant information to indicate why the label is True

False if both forms and narrative labels were False

narrative pages

- Input: String of text
- Tools: Tesseract, NER with SpaCy, DistilBERT
- Output: narrative label, score

Final classification for each document

Agenda

- **Project Background**
- **Data Hierarchy**
- **Our Pipeline**
 - Page Type Classification
 - Forms Pipeline
 - Narratives Pipeline
 - Training Data
 - Text Pre-processing
 - Model Justification + Selection
- **Conclusion and Future Work**

Strong results but avenues for future work remain

75%

less time required
for journalists to
review a case with
automation

Now journalists can better prioritize
which documents to look at first,
fast-tracking the annotation process

Areas for future work:

More granular classification

Can we classify at the page level, instead of document-level?

Multi-label classification

Does a document mention “51-50”? What about a specific substance?

Latency reduction in Azure

Currently takes 10 seconds to analyze a form page—moving documents onto Azure storage containers will reduce latency

Thank You! Questions?

Team BigLocal

Stanford DSSG

August 2022

Retrieve Google Drive files

- **Input:** County name
- **Output:** File path to each given county folder

Santa Clara County
Santa Cruz County
Shasta County
Sierra County

form pages

Page Type Classification

- **Input:** PDF file paths
- **Tools:** LayoutLMV2 Classifier
- **Output:** classification label per page (*form, narrative, interview, image, or other*)

rep_image	Label	doc	pg_num
./Shasta County/ima narrative	narrative	2022 Shasta County	16
./Shasta County/ima narrative	narrative	2022 Shasta County	48
./Shasta County/ima other	other	2022 Shasta County	15

narrative pages

Forms Pipeline

- **Input:** file path to form pages
- **Tools:** Azure Form Recognizer
- **Output:** form label, label justification

label	rel_info
False	NaN
True	{'signs of drug impairment': 'Page #1: [936.0,...
False	San Jose Police Department Property & Evidence...

narrative text from forms

Final classification for each document

Narratives Pipeline

- **Input:** String of text
- **Tools:** Tesseract, NER with SpaCy, DistilBERT
- **Output:** narrative label, score

Case Files (Raw Data) (%cd 'content/gdrive/MyDrive/Case Files (Raw Data)')

- dssg_biglocal
 - client_facing_pipeline.ipynb
 - model
 - mylib_dssg
 - form_recognizer.py
 - model.py
 - narrative_pipeline.py → for inference
 - narrative_training_pipeline.py → for training
 - label_agg.py
 - mylib (Hellina modified)
- Shasta County
 - 2022 Shasta County DA
 - Pdf1
 - Pdf2
 - Anderson PD
 - AndersonPDEricHaynes
 - **DIL report.pdf**
 - **CA News Coalition.pdf**
 - images
 - 2022 Shasta County DA__Pdf1-01.jpg
 - 2022 Shasta County DA__Pdf1-02.jpg
 - Anderson PD__AndersonPDEricHaynes__**DIL report.pdf**
 - Anderson PD__**CA News Coalition.pdf**