



Managing and Curating the AI Grand Challenge for Resilience

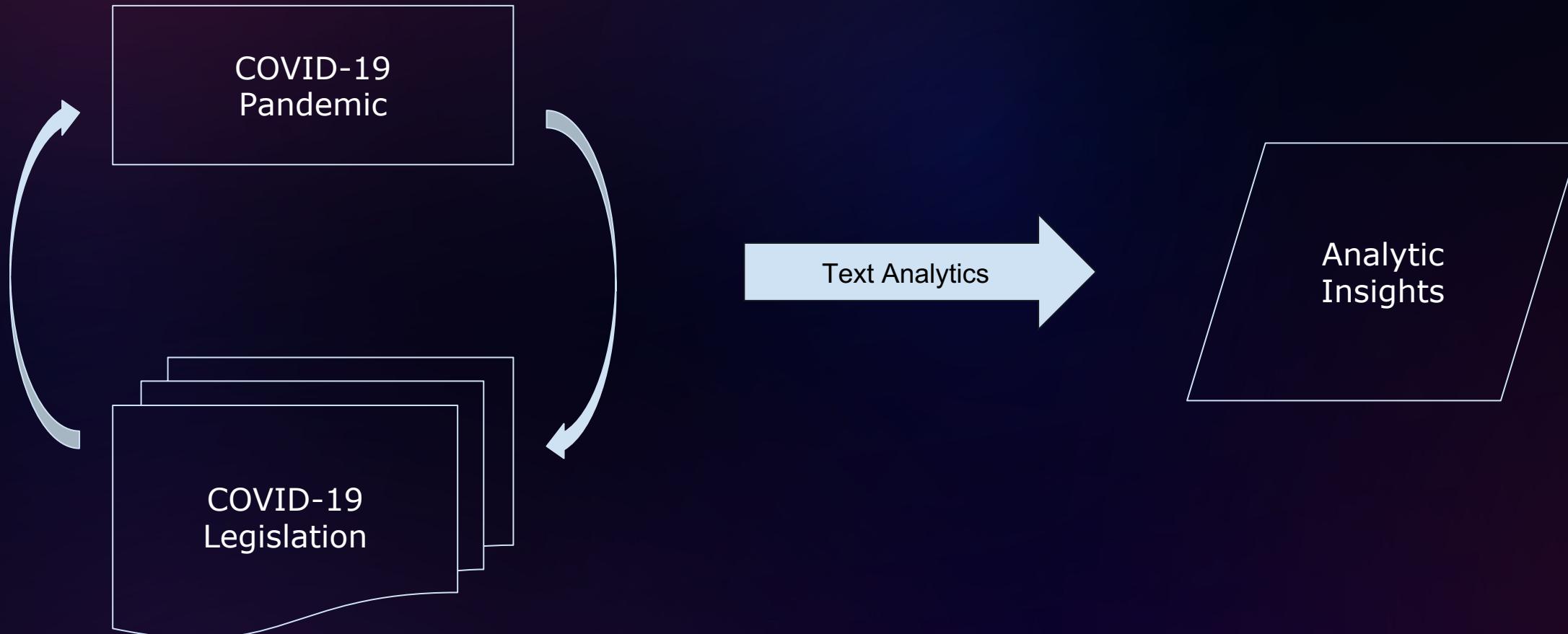
Bryan Lane, Maria Danilova, Steve
Babitch, Jarah Meador, Joseph Raeteno,
Melisa Marcovitz

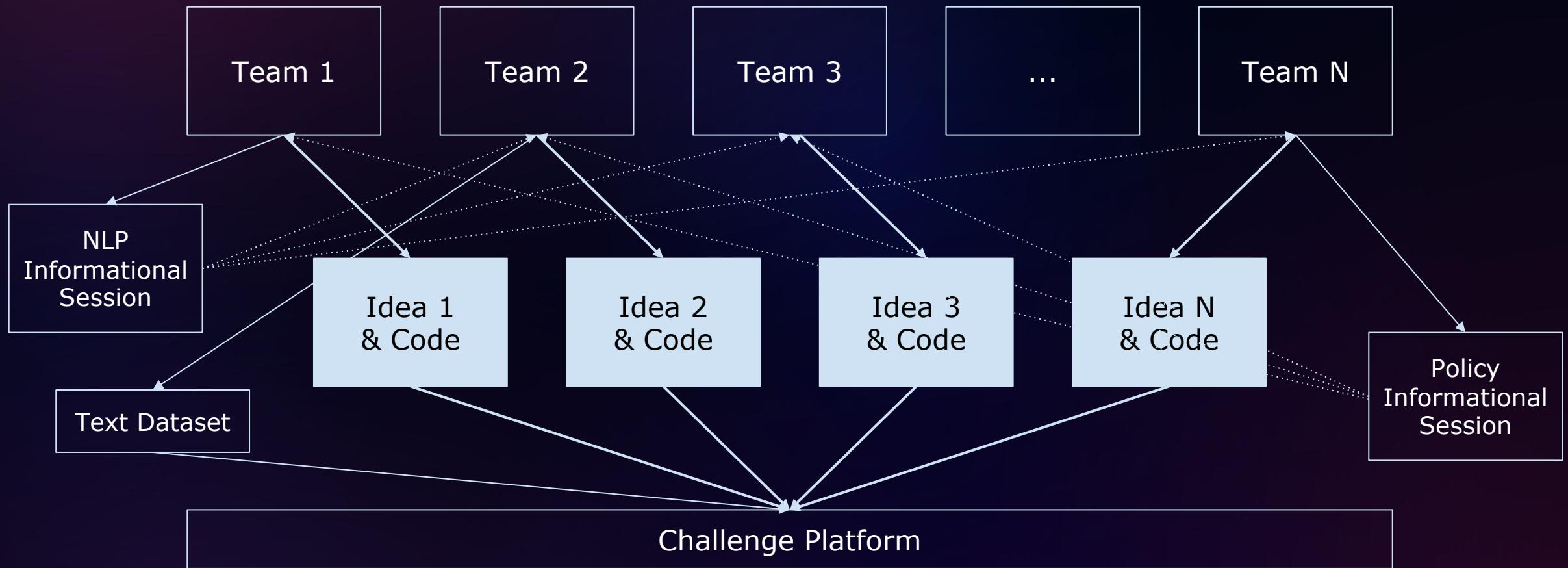
Presented on: 08/13/2021

Presented by:

Akhil Kondepudi -- University of Michigan, 2022

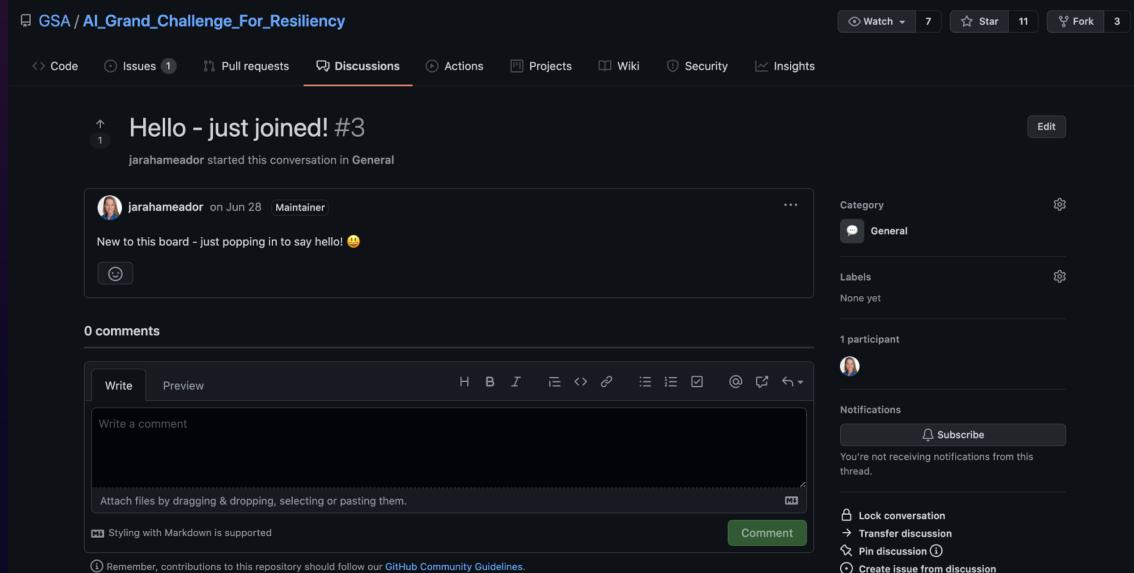
AI GRAND CHALLENGE FOR
RESILIENCE



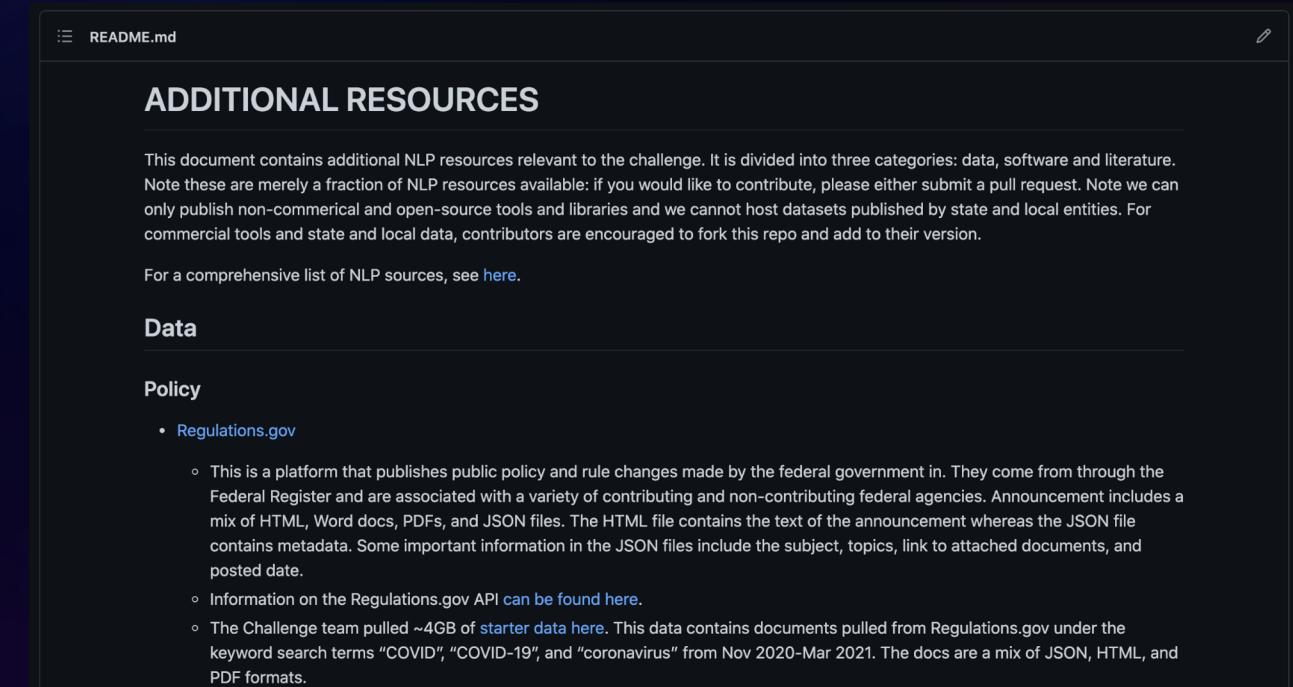


Hosting Platform: GitHub

AI GRAND CHALLENGE FOR
RESILIENCE



A screenshot of a GitHub discussion board titled "Hello - just joined! #3". The board shows a single message from user "jarahameador" on June 28, 2021, stating "New to this board - just popping in to say hello! 😊". The interface includes standard GitHub navigation links like Code, Issues, Pull requests, Discussions, Actions, Projects, Wiki, Security, and Insights. On the right side, there are settings for the discussion, including Category (General), Labels (None yet), Participants (1 participant), Notifications (Subscribe), and moderation options (Lock conversation, Transfer discussion, Pin discussion, Create issue from discussion).



A screenshot of a GitHub repository page for "GSA / AI_Grand_Challenge_For_Resiliency". The main content is a README.md file. It features a large heading "ADDITIONAL RESOURCES" followed by a block of text explaining the purpose of the document and encouraging contributions. Below this, there are sections for "Data" and "Policy". The "Policy" section includes a bullet point for "Regulations.gov" with several sub-points detailing its nature and the challenge team's use of its data.

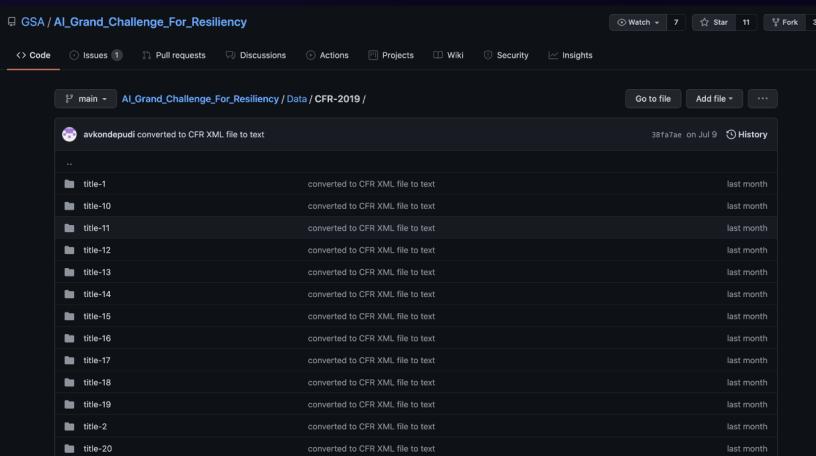
This document contains additional NLP resources relevant to the challenge. It is divided into three categories: data, software and literature. Note these are merely a fraction of NLP resources available: if you would like to contribute, please either submit a pull request. Note we can only publish non-commercial and open-source tools and libraries and we cannot host datasets published by state and local entities. For commercial tools and state and local data, contributors are encouraged to fork this repo and add to their version.

For a comprehensive list of NLP sources, see [here](#).

Data

Policy

- [Regulations.gov](#)
 - This is a platform that publishes public policy and rule changes made by the federal government. They come from through the Federal Register and are associated with a variety of contributing and non-contributing federal agencies. Announcement includes a mix of HTML, Word docs, PDFs, and JSON files. The HTML file contains the text of the announcement whereas the JSON file contains metadata. Some important information in the JSON files include the subject, topics, link to attached documents, and posted date.
 - Information on the [Regulations.gov API](#) can be found [here](#).
 - The Challenge team pulled ~4GB of starter data [here](#). This data contains documents pulled from [Regulations.gov](#) under the keyword search terms "COVID", "COVID-19", and "coronavirus" from Nov 2020-Mar 2021. The docs are a mix of JSON, HTML, and PDF formats.



A screenshot of a GitHub repository page for "GSA / AI_Grand_Challenge_For_Resiliency / Data / CFR-2019". The page lists a series of files named "title-1", "title-10", "title-11", "title-12", "title-13", "title-14", "title-15", "title-16", "title-17", "title-18", "title-19", "title-2", and "title-20". Each file entry includes a description: "converted to CFR XML file to text" and a timestamp: "last month". At the top of the list, there is a note: "avkondepot converted to CFR XML file to text" and a timestamp: "38fa7ae on Jul 9". The interface includes standard GitHub navigation links like Code, Issues, Pull requests, Discussions, Actions, Projects, Wiki, Security, and Insights.

Natural Language Processing

AI GRAND CHALLENGE FOR
RESILIENCE

Data

COVID-19 legislation, executive orders, titles

Methods

LDA, Word2Vec, deep learning

Software

spaCy, Stanza, NLTK

01

Data

Problem

Unstructured

Not “machine-readable”

Variability within class

Not all information is useful for ML purposes

Data Access

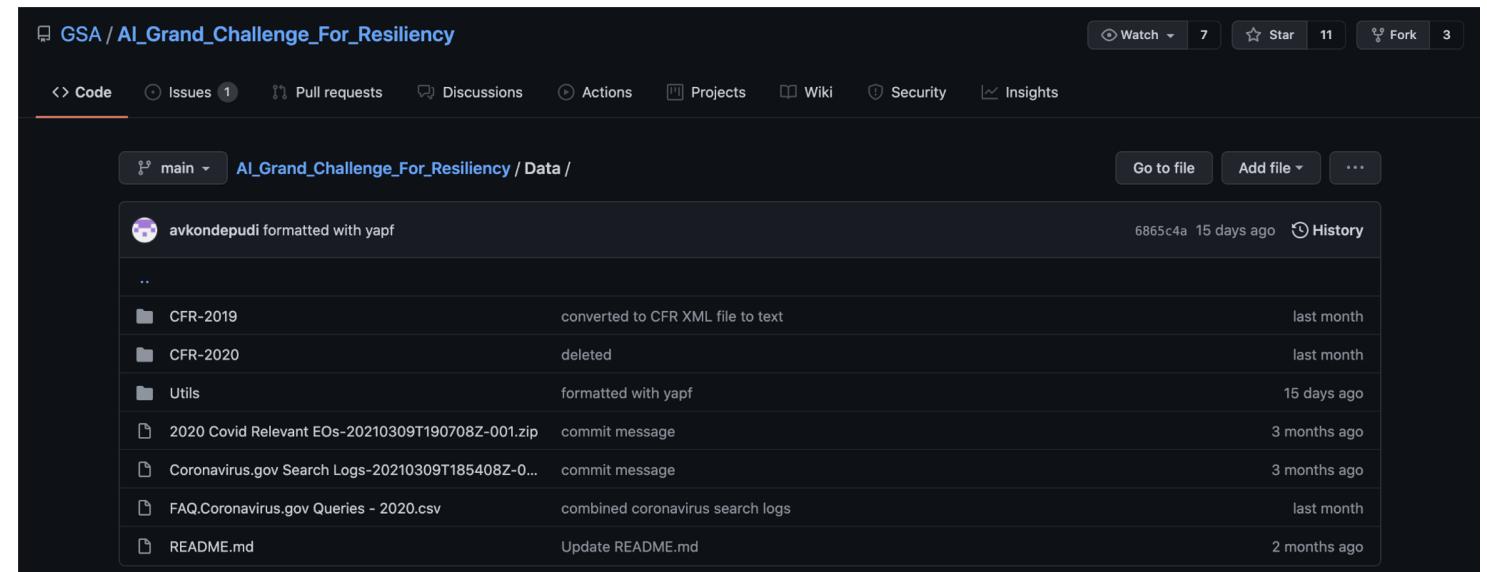
All data not present in one location

Exists across variety of formats (PDF, XML, etc.)

Hurdles to access & use data

Aggregated Data

Aggregated starter data for participants on GitHub



The screenshot shows a GitHub repository page for 'GSA / AI_Grand_Challenge_For_Resiliency'. The repository has 7 stars, 11 forks, and 3 issues. The main branch is 'main'. The 'Data' folder contains several files and subfolders, all committed by 'avkondepudi' using yapf. The commits are as follows:

Commit	Message	Date
avkondepudi formatted with yapf	converted to CFR XML file to text	last month
CFR-2019	deleted	last month
CFR-2020	formatted with yapf	15 days ago
Utils	commit message	3 months ago
2020 Covid Relevant EOss-20210309T190708Z-001.zip	commit message	3 months ago
Coronavirus.gov Search Logs-20210309T185408Z-0...	combined coronavirus search logs	last month
FAQ.Coronavirus.gov Queries - 2020.csv	Update README.md	2 months ago
README.md		

02

Methods

Problem

Natural language processing has had tremendous advances recently

However, applications to legislative data are mostly unexplored

- Methods may not be appropriate for use across multiple domains
- Results may require policy knowledge for interpretation

Literature Conferences

Compiled all resources and posted to GitHub page

A Frustratingly Easy Approach for Entity and Relation Extraction

Zexuan Zhong Danqi Chen
Department of Computer Science
Princeton University
{zzhong, danqic}@cs.princeton.edu

LEGAL-BERT: The Muppets straight out of Law School

Ilias Chalkidis †‡ Manos Fergadiotis †‡
Prodromos Malakasiotis †‡ Nikolaos Aletras * Ion Androultsopoulos †‡
† Department of Informatics, Athens University of Economics and Business
‡ Institute of Informatics & Telecommunications, NCSR “Demokritos”
* Computer Science Department, University of Sheffield, UK
[ihalk, fergadiotis, ruller, ion]@aueb.gr
n.aletras@sheffield.ac.uk



JURIX

The Foundation for Legal Knowledge Based Systems

Policy & NLP Experts

Helped curate list of experts to hold information sessions, give domain-specific advice to participants, etc.

Regulations.gov API



Office of Government-wide Policy Overview



Sample Submission

Provide example of implementation of grading rubric, ideas

Part I: Team Information & Overview

Team Information

Bryan Lane¹, bryan.lane@gsa.gov
Akhil Kondepudi¹, akhil.kondepudi@gsa.gov
Asitang Mishra², asitang.mishra@jpl.nasa.gov
Masha Danilova¹, maria.danilova@gsa.gov

¹Technology Transformation Services, General Services Administration

²Jet Propulsion Laboratory, NASA

Contributions

Akhil Kondepudi designed and wrote the submission. Mr. Bryan Lane and Mr. Asitang Mishra reviewed the submission. Ms. Masha Danilova was instrumental in facilitating collaborations for our project.

Statement of Interest

The team has a wide range of experience in both natural language processing and policy. This challenge provided us with an opportunity to test recent NLP advances to the policy domain and observe the results. Additionally, the team is extremely interested in the impact of regulations passed during the COVID-19 pandemic on the pandemic.

Part II: Submission Information

Title

Understanding regulations passed during the COVID-19 pandemic using text embeddings

Link to Forked Repository

<link>

Submission Overview

TTS AI Team's "Text Embedding" submission uses the Code for Federal Regulations, deep metric learning techniques like contrastive learning, and transformer models to characterize government-wide policy changes as a result of the COVID-19 pandemic.

Part III: Narrative Whitepaper

Project Overview

The COVID-19 pandemic spurred multiple regulations designed to protect industries' interests from the ensuing impact. These regulations can be found in the Code for Federal Regulations, codification of the general and permanent rules published in the Federal Register by the departments and agencies of the Federal Government produced by the Office of the Federal Register (OFR) and the Government Publishing Office. Representation learning in the policy domain has been sparse, so we address this bottleneck by training various models with specific pretexts on this dataset to create embeddings at both the document and sentence level. We then tested the efficacy of these modes by training a one-layer classifier on top of the frozen model. Finally, we plotted these embeddings under various classifications to better understand the effect of the pandemic on regulations.

Related Work

Metric Learning

The type of representation learning used in this submission is metric learning. This method aims to learn embeddings such that similar data are closer to each other in the embedding space, and vice versa. Current metric learning applications in the deep learning space first train deep learning models to extract rich feature representations using specific tasks known as "pretasks." These tasks are not relevant to the main purpose of the model. Next, the model is finetuned on a "downstream" task to test the quality of features learned, often by training a linear classifier on top of the frozen pretrained model.

Contrastive Learning

In recent years, contrastive learning has transformed the field of representation learning, especially in the field of computer vision. The idea of contrastive learning is that using an data point as an "anchor" makes similar, or "positive," data points closer to the anchor in an embedding space and different, or "negative," data points further away from the anchor. When doing unsupervised learning, an example of a positive sample may be a transformation of the anchor, whereas a negative sample would be a transformation of another sample in the anchor's mini-batch. This seemingly facile approach can be compounded upon by using the labels: now, positive samples are transformations of samples in the same class as the anchor and vice versa. This technique has led to significant improvements in the field of computer vision. Its application to NLP has been slow due to the discrete nature of the textual data, but preliminary results have already shown improvements in this space as well. For example, Giorgi et al. utilized a contrastive approach to unsupervised learning on text, which resulted in similar results as supervised and semi-supervised methods on downstream tasks.

The BERT Model

Specifics can be found in Delvin et al. (<https://arxiv.org/pdf/1810.04805.pdf>)

Applied Metric Learning

<MLM + NSP; Transformers>

03

Software

Current Landscape

Libraries

General, multipurpose NLP libraries

Tools

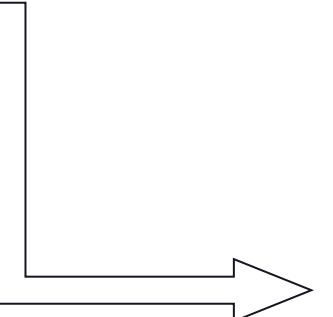
Perform specific task extremely well

Implementations

Open-source code for state-of-the-art deep learning models

Regulations Parser

```
<?xml version="1.0" encoding="UTF-8"?>
<CFRDOC xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance" xsi:noNamespaceSchemaLocation="CFRMergedXML.xsd">
<AMDDATE>May 2, 2018</AMDDATE>
<FMTR>
<TITLEPG>
<CODE/>
<PRTPAGE P="1"/>
<TITLENUM>Title 1</TITLENUM>
<SUBJECT>General Provisions</SUBJECT>
<REVISED>Revised as of January 1, 2019</REVISED>
<CONTAINS>Containing a codification of documents of general applicability and future effect</CONTAINS>
<DATE>As of January 1, 2019</DATE>
<PUB>
<P>Published by the Office of the Federal Register National Archives and Records Administration as a Special Edition of the Federal Register</P>
</PUB>
</TITLEPG>
<BTITLE>
<ONOTICE>
<PRTPAGE P="?ii"/>
<HD SOURCE="HED">U.S. GOVERNMENT OFFICIAL EDITION NOTICE</HD>
<HD SOURCE="HED">Legal Status and Use of Seals and Logos</HD>
<GPH DEEP="54" HTYPE="LEFT" SPAN="1">
<GID>archives.ai</GID>
</GPH>
<P>The seal of the National Archives and Records Administration (NARA) authenticates the Code of Federal Regulations (CFR) as the official codification of Federal regulations established under the Federal Register Act. Under the provisions of 44 U.S.C. 1507, the contents of the CFR, a special edition of the Federal Register, shall be judicially noticed. The CFR is prima facie evidence of the original documents published in the Federal Register (44 U.S.C. 1510).</P>
<P>It is prohibited to use NARA's official seal and the stylized Code of Federal Regulations logo on any republication of this material without the express, written permission of the Archivist of the United States or the Archivist's designee. Any person using NARA's official seals and logos in a manner inconsistent with the provisions of 36 CFR part 1200 is subject to the penalties specified in 18 U.S.C. 506, 701, and 1017.</P>
<HD SOURCE="HED">Use of ISBN Prefix</HD>
<P>This is the Official U.S. Government edition of this publication and is herein identified to certify its authenticity. Use of the 0-16 ISBN prefix is for U.S. Government Publishing Office Official Editions only. The Superintendent of Documents of the U.S. Government Publishing Office requests that any reprinted edition clearly be labeled as a copy of the authentic work with a new ISBN.</P>
</ONOTICE>
</BTITLE>
```



TITLE	DATE	TEXT	SECTNO	PART	PART DESC	PART NUMBER	SECT SUBJECT
Title 1—General Provisions	As of January 1, 2019	Any person may reproduce or republish, v	2.6	Pt. 2	PART 2—GENERAL INF	2	Unrestricted use.
Title 1—General Provisions	As of January 1, 2019	(a) Documents filed with the Office of the	3.2	Pt. 3	PART 3—SERVICES TO	3	Public inspection of docu
Title 1—General Provisions	As of January 1, 2019	(a) Pursuant to chapter 15 of title 44, Unit	5.1	Pt. 5	PART 5—GENERAL	5	Publication policy.
Title 1—General Provisions	As of January 1, 2019	The following documents are required to t	5.2	Pt. 5	PART 5—GENERAL	5	Documents required to be
Title 1—General Provisions	As of January 1, 2019	Whenever the Director of the Federal Reg	5.3	Pt. 5	PART 5—GENERAL	5	Publication of other docu
Title 1—General Provisions	As of January 1, 2019	The Federal Register serves as a daily su	5.5	Pt. 5	PART 5—GENERAL	5	Supplement to the Code.
Title 1—General Provisions	As of January 1, 2019	There shall be an edition of the Federal R	5.6	Pt. 5	PART 5—GENERAL	5	Daily publication.
Title 1—General Provisions	As of January 1, 2019	The Government Printing Office shall disti	5.7	Pt. 5	PART 5—GENERAL	5	Delivery and mailing.
Title 1—General Provisions	As of January 1, 2019	Without prejudice to any other form of citi	5.8	Pt. 5	PART 5—GENERAL	5	Form of citation.
Title 1—General Provisions	As of January 1, 2019	Each document published in the Federal I	5.9	Pt. 5	PART 5—GENERAL	5	Categories of documents
Title 1—General Provisions	As of January 1, 2019	Pursuant to section 1506 of title 44, Unit	5.1	Pt. 5	PART 5—GENERAL	5	Forms of publication.
Title 1—General Provisions	As of January 1, 2019	Each daily issue of the Federal Register s	6.1	Pt. 6	PART 6—INDEXES AND	6	Index to daily issues.
Title 1—General Provisions	As of January 1, 2019	Analytical subject indexes covering the cc	6.2	Pt. 6	PART 6—INDEXES AND	6	Analytical subject indexe
Title 1—General Provisions	As of January 1, 2019	(a) Each daily issue of the Federal Regis	6.3	Pt. 6	PART 6—INDEXES AND	6	Daily lists of parts affecte
Title 1—General Provisions	As of January 1, 2019	A monthly list of sections of the Code of F	6.4	Pt. 6	PART 6—INDEXES AND	6	Monthly list of sections af
Title 1—General Provisions	As of January 1, 2019	(a) The Director of the Federal Register m	6.5	Pt. 6	PART 6—INDEXES AND	6	Indexes, digests, and gu
Title 1—General Provisions	As of January 1, 2019	(a) Pursuant to chapter 15 of title 44, Unit	8.1	Pt. 8	PART 8—CODE OF FED	8	Policy.
Title 1—General Provisions	As of January 1, 2019	(a) Criteria. Each book of the Code shall t	8.3	Pt. 8	PART 8—CODE OF FED	8	Periodic updating.
Title 1—General Provisions	As of January 1, 2019	(a) The Director publishes a special editio	9.1	Pt. 9	PART 9—THE UNITED S	9	Publication required.
Title 1—General Provisions	As of January 1, 2019	(a) The Manual will contain appropriate in	9.2	Pt. 9	PART 9—THE UNITED S	9	Scope.
Title 1—General Provisions	As of January 1, 2019	The Director publishes a special edition o	10.1	Pt. 10	PART 10—PRESIDENTI	10	Publication required.
Title 1—General Provisions	As of January 1, 2019	The Director of the Federal Register shall	10.1	Pt. 10	PART 10—PRESIDENTI	10	Publication required.
Title 1—General Provisions	As of January 1, 2019	The Administrative Committee may autho	10.13	Pt. 10	PART 10—PRESIDENTI	10	Coverage of prior years.
Title 1—General Provisions	As of January 1, 2019	(a) The subscription price for the paper ec	11.2	Pt. 11	PART 11—SUBSCRIPTIC	11	Federal Register.
Title 1—General Provisions	As of January 1, 2019	(a) The online edition of the Manual, issu	11.4	Pt. 11	PART 11—SUBSCRIPTIC	11	The United States Gover
Title 1—General Provisions	As of January 1, 2019	The annual subscription price for the mon	11.7	Pt. 11	PART 11—SUBSCRIPTIC	11	Federal Register Index.

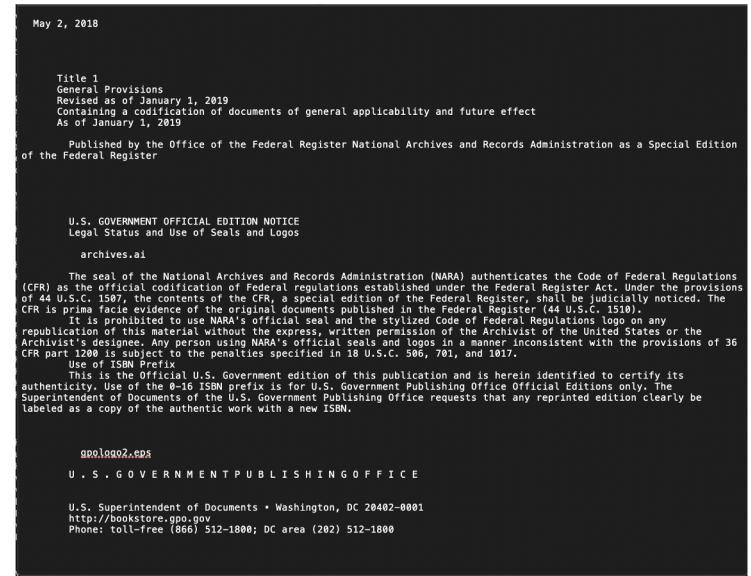
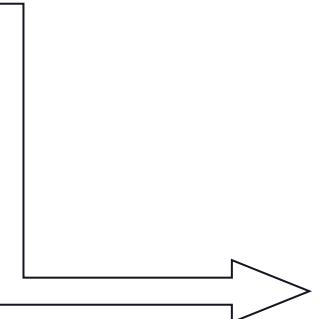
AI GRAND CHALLENGE FOR
RESILIENCE

Parses title data from
Code of Federal
Regulations

Regulations Parser

```
<?xml version="1.0" encoding="UTF-8"?>
<CFRDOC xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance" xsi:noNamespaceSchemaLocation="CFRMergedXML.xsd">
<AMMDATE>May 2, 2018</AMMDATE>
<FMTR>
<TITLEPG>
<CODE/>
<PRTPAGE P="1"/>
<TITLENUM>Title 1</TITLENUM>
<SUBJECT>General Provisions</SUBJECT>
<REVISED>Revised as of January 1, 2019</REVISED>
<CONTAINS>Containing a codification of documents of general applicability and future effect</CONTAINS>
<DATE>As of January 1, 2019</DATE>
<PUB>
<P>Published by the Office of the Federal Register National Archives and Records Administration as a Special Edition of the Federal Register</P>
</PUB>
</TITLEPG>
<BTITLE>
<ONOTICE>
<PRTPAGE P="?ii"/>
<HD SOURCE="HED">U.S. GOVERNMENT OFFICIAL EDITION NOTICE</HD>
<HD SOURCE="HED">Legal Status and Use of Seals and Logos</HD>
<GPH DEEP="54" HTYPE="LEFT" SPAN="1">
<GID>archives.ai</GID>
</GPH>
<P>The seal of the National Archives and Records Administration (NARA) authenticates the Code of Federal Regulations (CFR) as the official codification of Federal regulations established under the Federal Register Act. Under the provisions of 44 U.S.C. 1507, the contents of the CFR, a special edition of the Federal Register, shall be judicially noticed. The CFR is prima facie evidence of the original documents published in the Federal Register (44 U.S.C. 1510).</P>
<P>It is prohibited to use NARA's official seal and the stylized Code of Federal Regulations logo on any republication of this material without the express, written permission of the Archivist of the United States or the Archivist's designee. Any person using NARA's official seals and logos in a manner inconsistent with the provisions of 36 CFR part 1200 is subject to the penalties specified in 18 U.S.C. 506, 701, and 1017.</P>
<HD SOURCE="HED">Use of ISBN Prefix</HD>
<P>This is the Official U.S. Government edition of this publication and is herein identified to certify its authenticity. Use of the 0-16 ISBN prefix is for U.S. Government Publishing Office Official Editions only. The Superintendent of Documents of the U.S. Government Publishing Office requests that any reprinted edition clearly be labeled as a copy of the authentic work with a new ISBN.</P>
</ONOTICE>
</BTITLE>

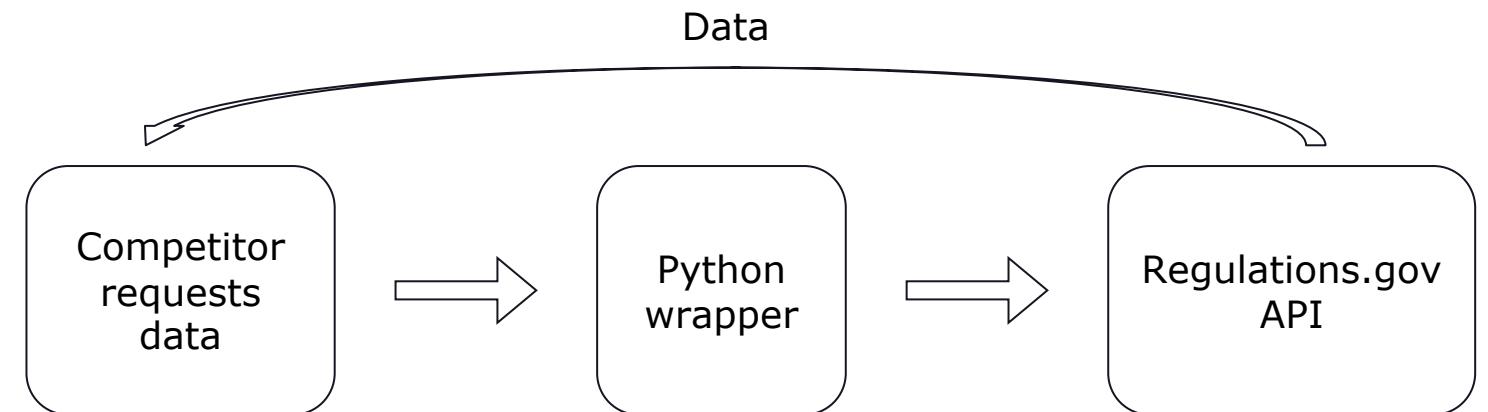
```



Regulations.gov API

Built Python wrapper around Regulations.gov API

Easier for participants to pull data, integrate into existing workflows



04

The Experience

Lessons I've Learned

AI GRAND CHALLENGE FOR
RESILIENCE

Government

Slow but necessary process

- Contractors
- Administration
- Checkpoints

Impact is undeniable, fulfilling

Work/life balance emphasized

Data Science

Drawing meaningful conclusions from results

Working alongside others (e.g., policy experts) to implement changes as a result of analyses

- Results from challenge will change how legislation occurs, lead to increased funding for specific initiatives, etc.

My work would not have been possible without the following people:

- The AIGC Team (Bryan Lane, Masha Danilova, and everyone else!)
- The CIF Team (Rachel Dodell, Ariana Soto, DJ Jain)
- TTS (Annabel Lombard, Molly McIntyre)

A big thank you to all of you!



AI GRAND CHALLENGE FOR
RESILIENCE

Questions?