# Agenda

1. Introductions
2. Hackathon Objective
3. Deliverables and Resources
4. General Information
5. Stages of Dashboard Implementation
6. Example Dashboard

https://hackhpc.github.io/ADMI22/

# Organizers

**Linda Hayden** - *ECSU/SGCI*
haydenl@mindspring.com

**Amy Cannon** - *Omnibond*
amycannon@omnibond.com

**Alex Nolte** - *University of Tartu*
alexander.nolte@ut.ee

**Boyd Wilson** - *Omnibond*
boyd@omnibond.com

**Je'aime Powell** - *TACC*
jpowell@tacc.utexas.edu

**John Holly** - *XSEDE*
jholly@sura.org
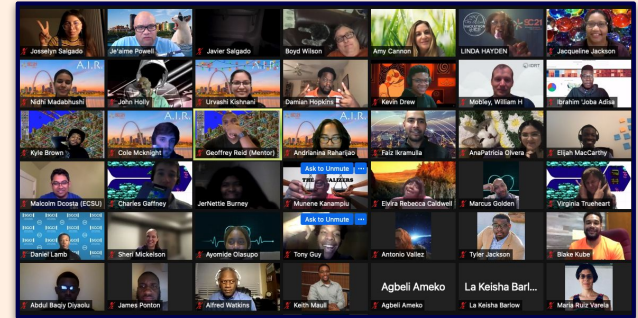
`https://hackhpc.github.io/ADMI22/`

# The Objective of HackHPC@ADMI

The hackathon aims to harness the resources, skills, and knowledge found in the HPC community in an effort to provide applied exposure towards students from 2-4 year post-secondary educational institutions. In short, the hackathon will provide HPC skills and training while targeting problems that directly affect the participants.

Develop knowledge about solutions to identified issues affecting them through application of data analysis/presentation or management.

## Student Outcomes

- Increased familiarity with data science in the cloud
- Experience collaborative software engineering
- Develop professional communication skills



`https://hackhpc.github.io/ADMI22/`

# Student Deliverables and Resources

## Deliverables:

- Source code Including Comments
- PDF of presentation
  - Team members with pictures
  - Use of HPC technology in the project
- Github Repository Link
  - README.md with project description

## Resources:

- Google Cloud (Provided Credits)
- Cloudy Cluster
- Most Commonly Used
  - Python
  - Jupyter Notebooks
  - Node.Js (JavaScript)
  - Repl.it (Collaborative Environment)
  - HTML
- Discord - https://discord.gg/ARg3vwWafF

`https://hackhpc.github.io/ADMI22/`

# General Information (the 3 T's)

- **Teams**
  - 4-5 Students
  - 1 Primary Mentor
  - 1 Technical Mentor
- **Time**
  - March 31st - April 4th
    - 3/31 @~7pm ET Event Start
      - "*The Draft*"
    - 4/[1-4] @ 11am ET & 7pm ET- Checkins
    - 4/4@6pm ET-Final Presentations

- **Topic Examples**
  - Data Analysis of COVID 19
  - Economic disparities and their effects on college participation
  - Genomics, Molecular Dynamics, or Weather Modeling in the Cloud.
  - Social Justice
  - AI-based Crowd Status
  - Public Data Management
  - Graduation Rates
  - Broadband Access
  - Insurance vs. Public Health Resilience

`https://hackhpc.github.io/ADMI22/`

# What is Data?

**data** noun, plural in form but singular or plural in construction, often attributive

🔖 Save Word

da·ta | \ ˈdā-tə 🔊 , ˈda- 🔊 *also* ˈdä- 🔊 \

### Definition of *data*

1  : factual information (such as measurements or statistics) used as a basis for reasoning, discussion, or calculation
   // the *data* is plentiful and easily available
   — H. A. Gleason, Jr.

   // comprehensive *data* on economic growth have been published
   — N. H. Jacoby

2  : information in digital form that can be transmitted or processed

3  : information output by a sensing device or organ that includes both useful and irrelevant or redundant information and must be processed to be meaningful

*Ref*: https://www.merriam-webster.com/dictionary/data

# OK, Yeah but... what is Data?

# OK, Yeah but... what is Data?

In the context of this presentation, data is information that you want to collect in a digital format for the purpose of analysis.

~ J. Powell

https://hackhpc.github.io/ADMI22/

# *I'm Scared to ask but...* What is a Dashboard?

A "*Dashboard*" frames a problem by telling a story using your data.

# The Dashboard Developmental Process

1. **Dashboard Design**
2. **Collect Data**
3. **Piping/Cleaning**
4. **Build Dashboard**
5. **Feedback/Test**

Design

Collect

Clean

Build

Test

https://hackhpc.github.io/ADMI22/

# Dashboard Design

- **Who** is the audience
- **What** information should they get from your dashboard / **What** question are you answering?
- **When** is the temporal connection between the dashboard and data [dynamic vs static data]
- **Where** - the platform?  (*desktop, server, kiosk*)
- **Why** - the goal for the whole project
  - Visualization - chart type
  - Pen and paper mockup

Design
Collect
Test
Clean
Build

**Output(s) from step:**
- Site mockup
- Clearly defined question(s)
- Platform to be used
- Type of data needed for analysis

https://hackhpc.github.io/ADMI22/

# Collect Data

*[Note: Third most time consuming process]*

- **Which datasets do you have access?**
- **What questions do you *WANT* to ask of the data?**
- **What questions *CAN* you answer from the available data you have?**
  - Alternate analysis/indirect correlations

Design
Collect
Test
Clean
Build

**Output(s) from step:**
- Dataset(s)
- Data Dictionaries
- Suggested analysis/correlation methods
- Dataset Documentation
- Database/Storage location(s)

https://hackhpc.github.io/ADMI22/

# Data Piping/Cleaning

**[*Note: Most time consuming process!*]**

- **Take raw data in**
- **Write scripts for necessary data transformations**
  - *Python, R, JupyterNotebook*
- **Identify data storage locations**
- **Handle moving data between locations**
- **Consider: data that changes over time**

Design
Collect
Clean
Build
Test

**Output(s) from step:**
- Clean Dataset(s)
- scripts for transformation
- output files
- database connections

# Data Piping/Cleaning - GIGO

*The reason this is the  MOST time consuming process!*

**GIGO = Garbage In, Garbage Out**

If your data is not properly organized and

"transformed" the results will likely not make sense!

- Data Validation

- Proper/Non-Repeating Headers

- Proper databases
  - Georeference-enabled

Design

Collect

Clean

Test

Build



THE THING ABOUT DATA IS...

GARBAGE IN,  GARBAGE OUT

`https://hackhpc.github.io/ADMI22/`

# Build Dashboard

*[Note: Second most time consuming process]*

- **Load outputs of data pipes/sources**
- **Code chart elements on page**
- **Code User interactivity**
  - **Data filters**
  - **Selection methods**
  - **Changing elements**

Design
Collect
Test
Clean
Build

**Output(s) from step:**
- Code used to build the dashboard
- Deployed dashboard locally or to a cloud service

https://hackhpc.github.io/ADMI22/

# Feedback/Testing

- **Demonstration to Client / Users**
  - Ideally a live deployed version
  - Screenshots / PDF better than nothing
- **Collect and Integrate feedback into next iterative development cycle**

**DID DASHBOARD TELL THE STORY / AID THE DECISION / ANSWER THE QUESTION?**

Design
Collect
Test
Clean
Build

**Output(s) from step:**
- Documented feedback
- Informed tasking for the next iteration(s) of the design

https://hackhpc.github.io/ADMI22/

# Iterate the process until done!



it's a never ending story



Design

Collect

Dashboard Development Process

Test

Clean

Build

# Dashboard Example

## Demo Time!!



**Example GitHub Repo:**
https://github.com/mepearson/texas_congress

**Deployed Heroku App:**
https://texas-congress.herokuapp.com/

`https://hackhpc.github.io/ADMI22/`

# Dashboard Design

**WHAT:**

- Dashboard to link TX residents with information for their US Congressional District

**DESIRED ELEMENTS:**

- Selectable map of Congressional Districts
- Display section for information related to selected District
- Table of clickable links to access District Information Files



https://hackhpc.github.io/ADMI22/

# Collect Data

Design

**Collect**

Test

Clean

Build



Get links to Congressional District maps from redistricting.capitol.texas.gov site

`https://hackhpc.github.io/ADMI22/`

# Data Piping/Cleaning

Design
Collect
Clean
Test
Build

### Texas Congress Website: Code development

**Python Libraries**

```
In [19]:  # Data Proessing
          import pandas as pd
          import geopandas as gpd

          # Data Visualization
          import plotly.express as px
```

**Texas Congressional District Map**

**ETL for Congress geospatial Data**

```
In [1]:   # Data file with geojson of US Congressional Districts
          # Congressional geospatial data downloaded from 2021 census.gov shapefiles
          # https://www.census.gov/cgi-bin/geo/shapefiles/index.php?year=2021&layergroup=C

          congress = 'C:/Users/lissa/Box/TACC/tx_congress/data/tl_2021_us_cd116.json'
```

```
In [2]:   # Process US Congressional geojson to extract TX data sand save TX only geojson

          gdf = gpd.read_file(congress)
          texas = gdf[gdf.STATEFP == '48']
          texas.reset_index(inplace=True)
          texas.to_file('texas_congress.geojson', driver='GeoJSON')
```

**Generate Data Files**

**CSV of Congressiona District and Redistricting Map pdf**

```python
## Create link to district map
district_map_link_prefix = 'https://wrm.capitol.texas.gov/fyiwebdocs/PDF/congres
district_map_link_suffix = '/m1.pdf'

cds = []
district_map_urls = []

for i in range(1,37):
    cd = str(i)
    cd_url = ''.join([district_map_link_prefix,str(i),district_map_link_suffix])
    if len(cd) == 1:
        cd = '0' + cd
    cds.append(cd)
    district_map_urls.append(cd_url)

district_dict={'CD116FP' : cds,
    'district_map_url' : district_map_urls,
    'type' : 'map',
    'filetype' : '.pdf',
    'description' : 'District Map from https://redistricting.capitol.texas.gov'
    }
district_files = pd.DataFrame(district_dict)

# Export data frame to csv
district_files.to_csv('district_files.csv')
```
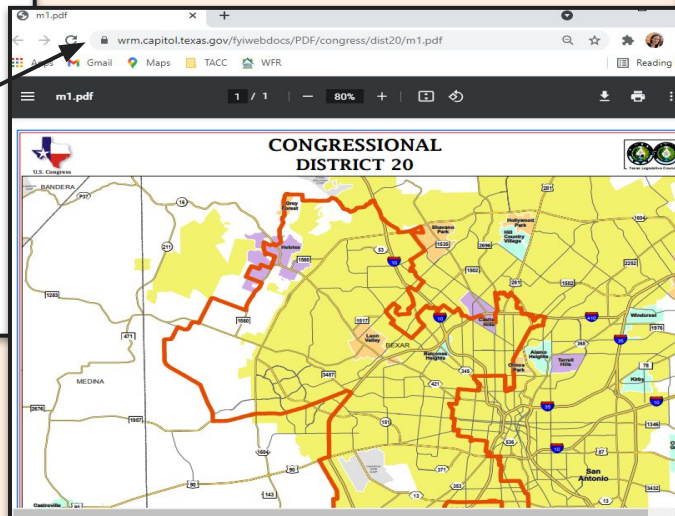
Use geopandas package in Jupyter notebook to extract Texas-only geojson

*Jupyter Notebook file available in assets folder of Github repo*

`https://hackhpc.github.io/ADMI22/`

# Build Dashboard

- **Write Dash code in IDE of choice**
- **Parts of App.py File:**
  - Python libraries
  - DATA Loading and DATA Visualizations
  - APP Layout – layout elements of page, similar to html
  - Callbacks – provide user interactivity / communication between elements
  - Run App

Design

Collect

Test

Clean

Build

```
109   # --------------------------------------------------
110   # CALLBACKS
111   # --------------------------------------------------
112
113   @callback(
114       Output('div-map-select', 'children'),
115       Output('div-files','children'),
116       Input('graph-map', 'clickData'))
117   def update_figure(clickData):
118       # Data for table of files
119       table_data_cols = ['Congress','State','District', 'File']
120       table_data = district_files[table_data_cols]
121
122       if clickData is None:
123           div_map = html.P('Select a Congressional district from the map at left to load the District Map')
124
125
126       # if District selected in map, display specialty map and filter files list
127       else:
128           # get value of district selected
129           cd = clickData['points'][0]['customdata'][0]
130           if cd[0] == '0': # remove leading 0
131               cd = cd[1:]
132           # get link to District map for selected district
133           cd_link = ''.join([district_map_link_prefix,cd,district_map_link_suffix])
134           div_map = html.Embed(src=cd_link,width="600",height="600",type="application/pdf")
135           # filter files table to district
```
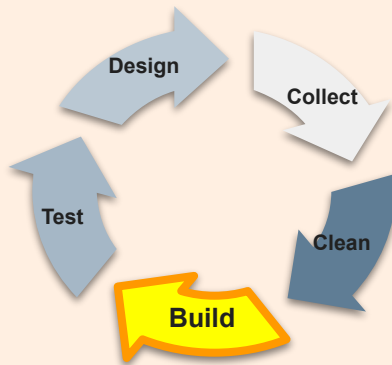
```
76    # --------------------------------------------------
77    # APP Layout
78    # --------------------------------------------------
79
80    external_stylesheets = [dbc.themes.LITERA]
81
82    app = Dash(__name__, external_stylesheets=external_stylesheets)
83
84    app.layout = html.Div([
85        dbc.Row([
86            html.H2('Texas Congressional District Information'),
87        ]),
88        dbc.Row([
89            dbc.Col([
90                dcc.Graph(
91                    id='graph-map',
92                    figure=map_fig,
93
94                ),
95            ],width=4),
96            dbc.Col([
97                html.Div(id='div-map-select'),
98                html.Div('Maps from https://redistricting.capitol.texas.gov/')
99            ],width=8),
100       ]),
101       dbc.Row([
102           dbc.Col([
103               html.Div(id='div-files'),
104           ])
105       ])
106   ])
107
```

`https://hackhpc.github.io/ADMI22/`

# Feedback/Testing



**Design** → **Collect** → **Clean** → **Build** → **Test**

**Deployed Heroku App:**
https://texas-congress.herokuapp.com/

`https://hackhpc.github.io/ADMI22/`

# Additional References

**Data Management**

- R for Data Science. Code in R / concepts useful any language
  [Welcome | R for Data Science (had.co.nz)](Welcome | R for Data Science (had.co.nz))

- Blog Overview (easy read): [Tidy data for efficiency, reproducibility, and collaboration (openscapes.org)](Tidy data for efficiency, reproducibility, and collaboration (openscapes.org))

- Original paper by Hadley Wickham (founder of R) who pioneered the concept of tidy data:
  - Official Paper: [Tidy data (had.co.nz)](Tidy data (had.co.nz))
  - informal and example code heavy (in R) version: [Tidy data • tidyr (tidyverse.org)](Tidy data • tidyr (tidyverse.org))

**Data Visualization**

- Chart Chooser — Juice Analytics - [https://www.juiceanalytics.com/chartchooser](https://www.juiceanalytics.com/chartchooser)

- Plotly graphing library - [https://plotly.com/python/](https://plotly.com/python/)

**Dash App**

- Dash App documentation - [https://dash.plotly.com/](https://dash.plotly.com/)

- Deploy to Heroku
  - integration from github [[https://devcenter.heroku.com/articles/github-integration](https://devcenter.heroku.com/articles/github-integration)]
  - Dash guidance / command line (scroll past Enterprise information to Heroku / free section) - [https://dash.plotly.com/deployment](https://dash.plotly.com/deployment)

# Questions and Concerns

Next Training Session:

## - Google / CloudyCluster - [6/26/22]

Schedule:

**https://hackhpc.github.io/ADMI22/schedule.html**

Presenters Contact Information:

**Je'aime Powell (*TACC*) - jpowell@tacc.utexas.edu**

`https://hackhpc.github.io/ADMI22/`