# Agenda

- **Introductions**
- **Hackathon Objective**
- **Deliverables and Resources**
- **General Information**
- **Data to Dashboard**

# Organizers

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

Boyd Wilson - *Omnibond*
boyd@omnibond.com

Amy Cannon - *Omnibond*
amycannon@omnibond.com

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

Linda Hayden - *ECSU*
haydenl@mindspring.com

HPC in the City: *St. Louis* HACKATHON

http://hackhpc.org/hpcinthecity

SC21
St.Louis, MO science & beyond.

# The Objective of HPC in the City: St. Louis

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 St. Louis 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

HPC in the City: St. Louis HACKATHON

http://hackhpc.org/hpcinthecity

SC21
St.Louis, science MO & beyond.

# Student Deliverables and Resources

**Deliverables:**

- **Source code Including Comments**
- **PDF of presentation**
  - Team members with pictures
  - Use of HPC technology in the project
  - Regional (St. Louis) implications of 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.com/invite/rSXasYKDwE**



Join the HPCHack
Discord using this
QR Code!

HPC
in the
City: St. Louis  **HACKATHON**

http://hackhpc.org/hpcinthecity

SC21
St.Louis, science
MO & beyond.

# General Information (the 3 T's)

- **Teams**
  - 4-5 Students
  - 1 Primary Mentor
  - 1 Specialist/Staff
- **Time**
  - November 4th - 8th
    - 11/4@~6pm ET Event Start
      - Team formation
    - 11/[5-8] @ 11 ET & 6pm ET- Checkins
    - 11/8@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

HPC in the City: St. Louis HACKATHON

http://hackhpc.org/hpcinthecity

SC21
St.Louis, MO science & beyond.

*Presenter:* Melissa Pearson (TACC)

# Creating a Data-based Dashboard Application

HPC in the City: St. Louis HACKATHON

http://hackhpc.org/hpcinthecity

SC21
St.Louis, science
MO & beyond.

# Talk Structure

| | |
|---|---|
| **Data to Dashboard Workflows** | • Workflow: Perception....<br>• .... vs. Reality<br>• Step by Step walkthrough |
| **Example Dash App** | • Simple application use case: TX Congressional District Info<br>• Example of workflow and products at each step<br>• Quick overview of deploying Github repo to Heroku |
| **Data Piping: In Detail** | • Data Cleaning and Wrangling<br>• E(xtraction), T(ransformation) and L(oad)<br>• "Tidy" Data |

# Data to Dashboards: Workflows

# How people imagine the workflow....

Get Data

- I have my inputs!

Make Dashboard

- Just wire up the code!

# ....and data visualization development....

**Get Data**
- I have my inputs!

**Clean and Prep data**
- Shouldn't be too hard..

**Data Visualizations**
- How many options can there be?

**Make Dashboard**
- Just wire up the code!

# ....and scale it to work effort....

*Expect data wrangling = 80% of your project.*

# Development Cycle



- Agile development is *iterative*.

- Quick cycling / rapid feedback.

- Keeps project from veering too far off course before a chance for correction.

# Development Cycle

**STEPS**

Dashboard Design

Collect Data

Data Piping

Build Dashboard

User Feedback

# Development Cycle



STEPS

1. Basic design / purpose

# Development Cycle



**STEPS**

1. Basic design / purpose
2. **Get your data**

# Development Cycle



**STEPS**

1. Basic design / purpose
2. Get your data
3. **Move, clean and shape the data.**

# Development Cycle



**STEPS**

1. Basic design / purpose
2. Get your data
3. Move, clean and shape the data.
4. **Build and deploy data visualization dashboard**

# Development Cycle



**STEPS**

1. Basic design / purpose
2. Get your data
3. Move, clean and shape the data.
4. Build and deploy data visualization dashboard
5. **Demo with User for feedback.**

*...and Repeat*

# Data to Dashboards: Workflow Steps - *Detailed*

# Dashboard Design



- A Dashboard frames the problem and tells the story in your data

  - Who – is the audience

  - What – information should they get from your dashboard

  - When – temporal connection between the dash and data [live vs not]

  - Where – platform, desktop vs. mobile

  - Why – goal for the whole project

- Data Visualization: the right chart for the need (see Resources)

## *OUTPUT FROM STEP: Site mockup*

WHITEBOARD
PEN AND PAPER

# Collect Data



- What data do you have access to

- What questions do you *want* to ask of the data?

- What questions *can* you answer from the data you have?

- *"You can't always get what you want"*

**OUTPUT FROM STEP: Documentation, data dictionaries, file directory**

| Github Google Drive Local File Directory |
|---|

# Data Piping



- Take raw data in

- Write scripts for necessary data transformations

- Identify data storage locations

- Handle moving data between locations

- *Consider: data that changes over time*

**OUTPUT FROM STEP: scripts for transformation, output files, database connections**



Github
Conda
Jupyter Notebook

# Build Dashboard

- Load outputs of data pipes into app

- Layout elements on page

- Wire-up User interactivity
  - Data filters
  - Selections
  - Changing elements

**OUTPUT FROM STEP: Code to build dashboard, deployed locally or to the cloud**

Dashboard Design

Collect Data

User Demo

BUILD DASHBOARD

Data Piping

*Form follows Function*

Dash
IDE [Atom]
Github

# User Feedback

- Demonstrate to Client / Users
  - Ideally deployed version that they can use
  - Try to WATCH user while they use system
  - Screenshots / PDF better than nothing
  - Collect feedback

- Integrate feedback into next iterative development cycle

- **KEY QUESTION**: DID DASHBOARD TELL THE STORY / AID THE DECISION

*OUTPUT FROM STEP: Documented feedback to inform next iterations of design*

Dashboard Design

Collect Data

USER FEEDBACK

Data Piping

Build Dashboard

Website [Heroku]
Powerpoint
PDF

# Data to Dashboards: Dash App Example

# 1. Dashboard Design

WHAT:

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

DESIRED ELEMENTS:

1. Selectable map of Congressional Districts

2. Display section for information related to selected District

3. Table of clickable links to access District Information Files

# 1. Dashboard Design: mock-up
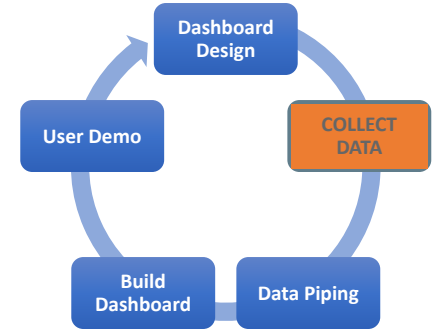


1. Selectable Map

2. Table of District Files with linked URLs
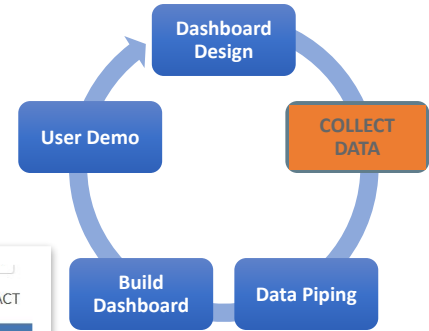
3. District Specific Information

# 2. COLLECT DATA

Identify online Data repositories.

- Geospatial data – Congressional District Geometries (map)

- PDF maps of individual districts

- Information on individual members
  *[data source identified, not yet collected]*

# 2. COLLECT DATA: Geospatial data

# 2. COLLECT DATA: District files



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

# 2. COLLECT DATA: member info *(to do)*



House members website

Examine content with developer tools to determine how to scrape site

# 3. DATA PIPING



- For data that needs transformation: move data where needed, transform, save new data output

- Load data files in local data directory

- Jupyter notebook to perform data transformations and develop data visualizations

# 3. DATA PIPING: geospatial



Import US district shapefile from Census to Mapshaper.org
☐ geojson

Use geopandas package in Jupyter notebook to extract Texas-only geojson

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

# 3. DATA PIPING: district maps file



1. Determine pattern to generate map url for each district

2. Generate dataframe of congressional district and links to file urls.

3. Export dataframe to csv file

## 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')
```

# 4. BUILD DASHBOARD

- Write Dash code in IDE of choice

- Parts of App.py File:

    1. Python libraries

    2. DATA Loading and DATA Visualizations

    3. APP Layout – layout elements of page, similar to html

    4. Callbacks – provide user interactivity / communication between elements

    5. Run App

# 4. BUILD DASHBOARD



## Develop Locally

1. Clone Github Repo [mepearson/texas_congress]
2. Create and launch virtual environment
3. Develop changes locally using IDE of choice [atom, vstudio, etc.]
4. Command Line: > python app.py

## Deploy to Heroku

1. Push changes to Github

2. Create site on Heroku linked to Github repo

3. Manually Deploy

# 4. BUILD DASHBOARD: layout & callbacks



```
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.
99           ],width=8),
100      ]),
101      dbc.Row([
102          dbc.Col([
103              html.Div(id='div-files'),
104          ])
105      ])
106  ])
107
```

```
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
```
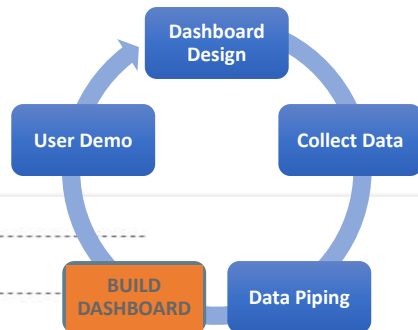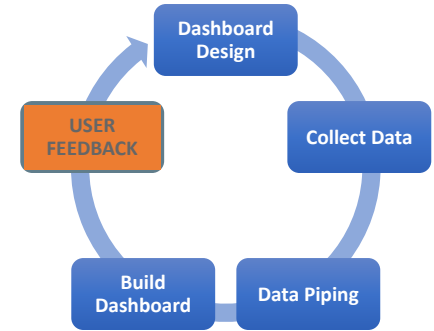
# 5. USER FEEDBACK

- Formal vs. Informal

- Decide how will collect feedback

- Set a process to track and prioritize feedback [spreadsheet, github issues, etc.]

- Use this to assign tasks

- Product Management is key. Tasks should have clear decision makers / authority for action.

# Data to Dashboards: DATA PIPING

# Data Cleaning and Wrangling

- "You can't always get what you need"

- Data arrives in many – usually Messy! - forms

- Data dashboards require standardized data

- Clean and wrangle what you get --> what you need

- What you need:
  - Quality controlled data...
  - Where you can access it..
  - In the right structure.          *Garbage in --> Garbage out*

# Quality Control Data

- If collecting your own data: make sure it's good.

- Pre-existing data:
  - review for quality [assess provider, etc.]
  - handle missing values

- Cite your data. Track sources and date acquired.

# Move and Shape your data [ETL]

- ETL has specific meaning, but abstract concept useful here

- ETL  = Extract, Transform, Load

- Extract data from data sources [websites, databases, files]

- Transform data into the necessary format
    - Clean – handle null values and bad data
    - Shape – wrangle data into required structure

- Load data into intermediate repository [database, file] or directly into dashboard

# Shaping data: make it 'Tidy'

- Everything assumes clean structured data, usually in 'Tidy' format

- Concepts language agnostic – learn on what you know best

- See 'Resources' slides for in depth discussion and references

# TIDY Data: what it is

Each variable is a column.

Each observation is a row.

Each value is a cell.



variables

observations

values

# Shaping data: how to make it 'Tidy'

How to Tidy data: tools available

- R: Tidyverse packages

- Python: pandas and associated packages

- Excel: pivot tables and other built in functions to add text

# Data to Dashboards: Resources

**Data Management**

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

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

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

**Data Visualization**

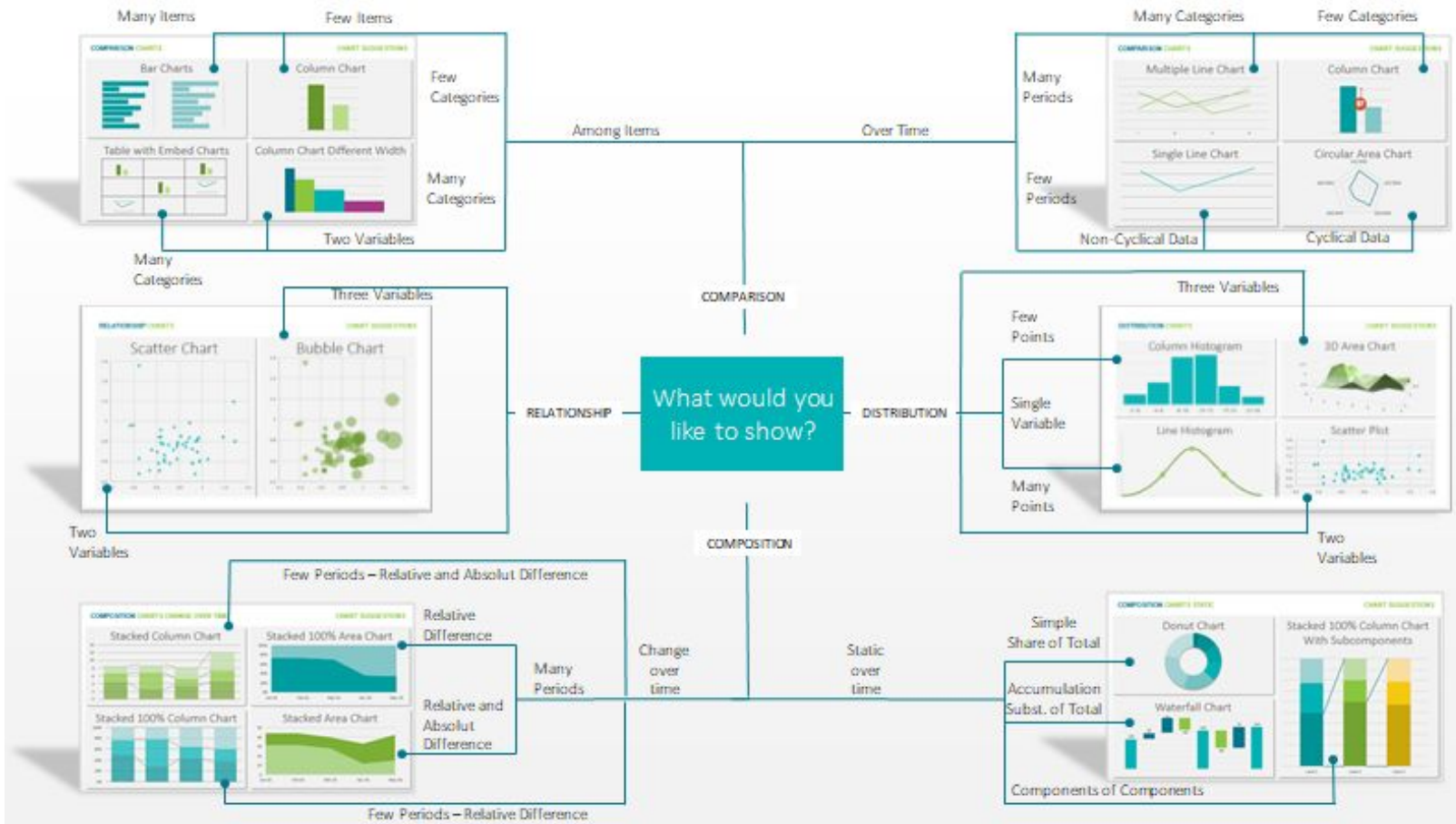- Chart Chooser — Juice Analytics - [https://www.juiceanalytics.com/chartchooser](https://www.juiceanalytics.com/chartchooser)
- Plotly graphing library - [https://plotly.com/python/](https://plotly.com/python/)

**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)

# Dashboard Design: Chart Selection

# Questions and Concerns

Next Training Sessions:

- **Beginning to End Project Example -** [*10/28/21*]

- **Kick-Off** - [*11/4/21*]

Schedule:
**https://jeaimehp.github.io/HackHPC-HPCintheCity21/**

Presenter Contact Information:

**Melissa Pearson (TACC) -** Lissa.Pearson@austin.utexas.edu



HPC in the City: *St. Louis* **HACKATHON**

**http://hackhpc.org/hpcinthecity**

SC21
St.Louis, MO | science & beyond.