# COMPUTING WITH CORGIS:
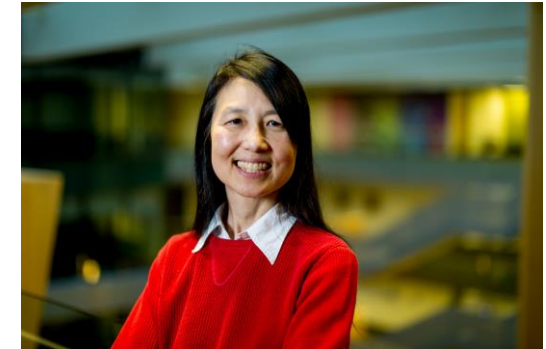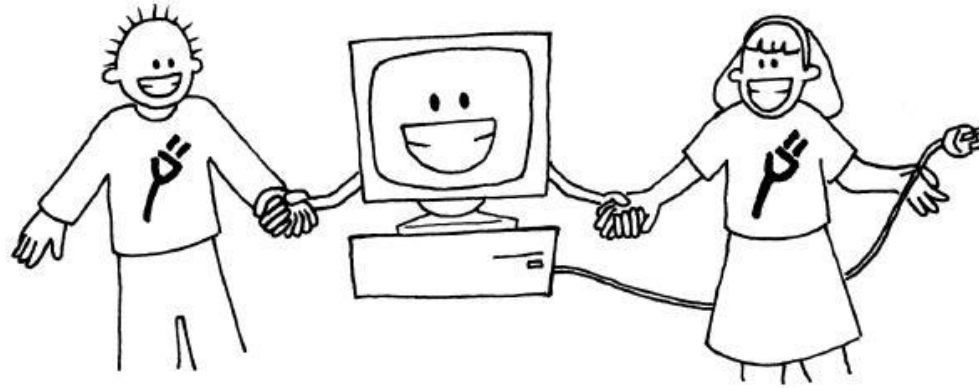# DIVERSE, REAL-WORLD DATASETS FOR INTRODUCTORY COMPUTING

Austin Cory Bart, Ryan Whitcomb,
Dennis Kafura, Clifford A. Shaffer, Eli Tilevich
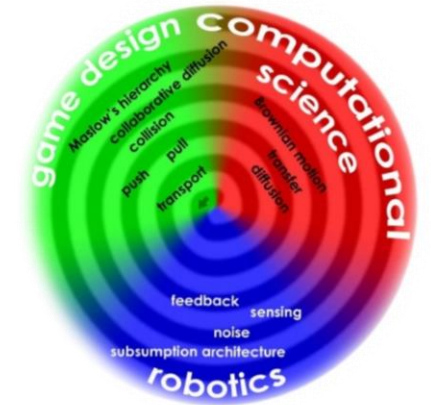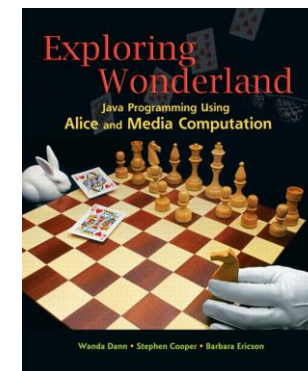
Virginia Tech

# Overview

- Bringing real-world data into introductory computing classes

- Via a new system that manages and produces datasets

- In order to motivate non-computing majors
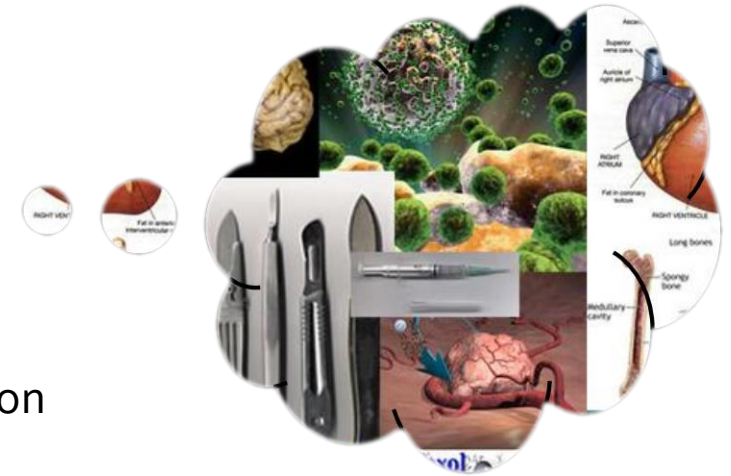
# Computer Science For All

# Diverse Majors                    … with Rich Knowledge



Animal Sciences

English

Biological Sciences

Education

Chemistry

Theater Arts

Building Construction

History

# (1) No Prior Background



"I've never done this before."

# (2) Low Self-efficacy



"I have no idea how to do this!"

# (3) Unclear on Why



"Why am I doing this?"
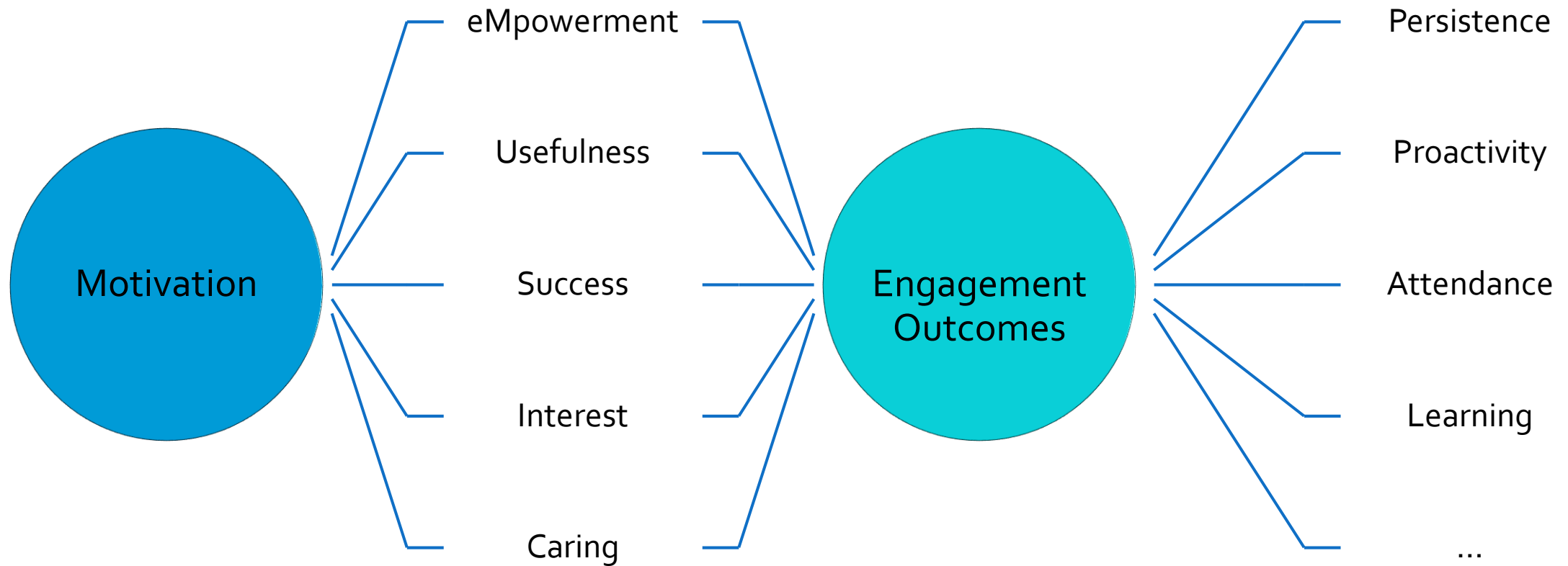
# MUSIC Model of Academic Motivation

Students are more motivated when they **perceive** that:

1. they are **eMpowered**,

2. the content is **Useful** to their goals,

3. they can be **Successful**,

4. they are **Interested**, and

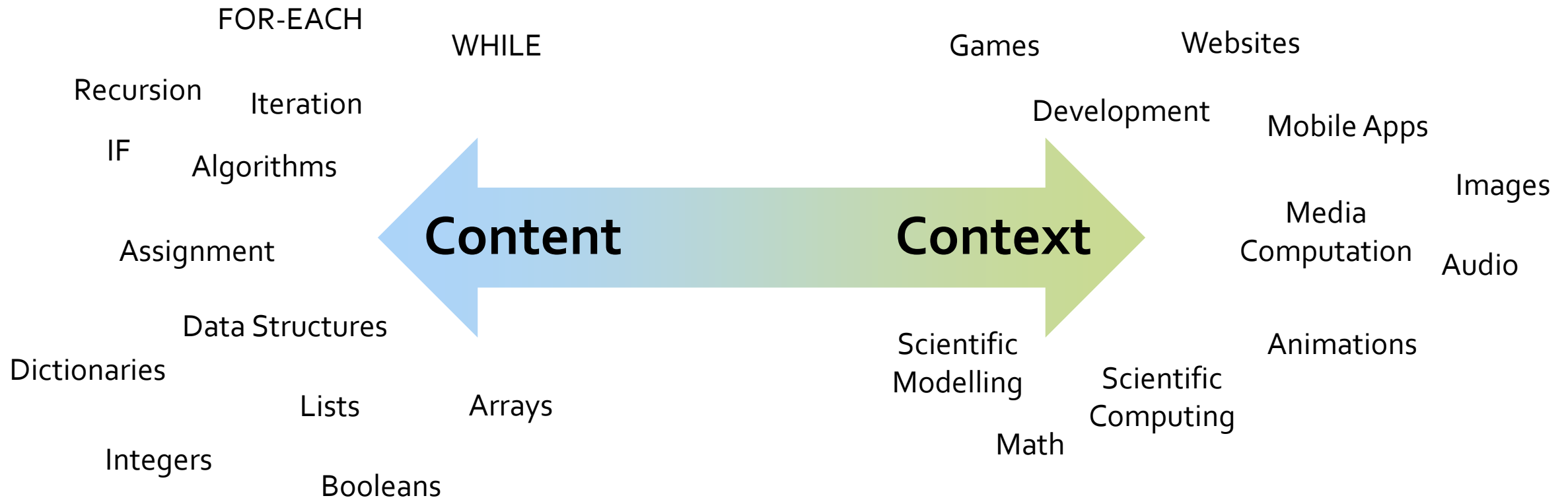5. they feel **Cared** for by others in the learning environment

*B. D. Jones. Motivating students to engage in learning: The MUSIC model of academic motivation. International Journal of Teaching and Learning in Higher Education, 21(2):272–285, 2009.*
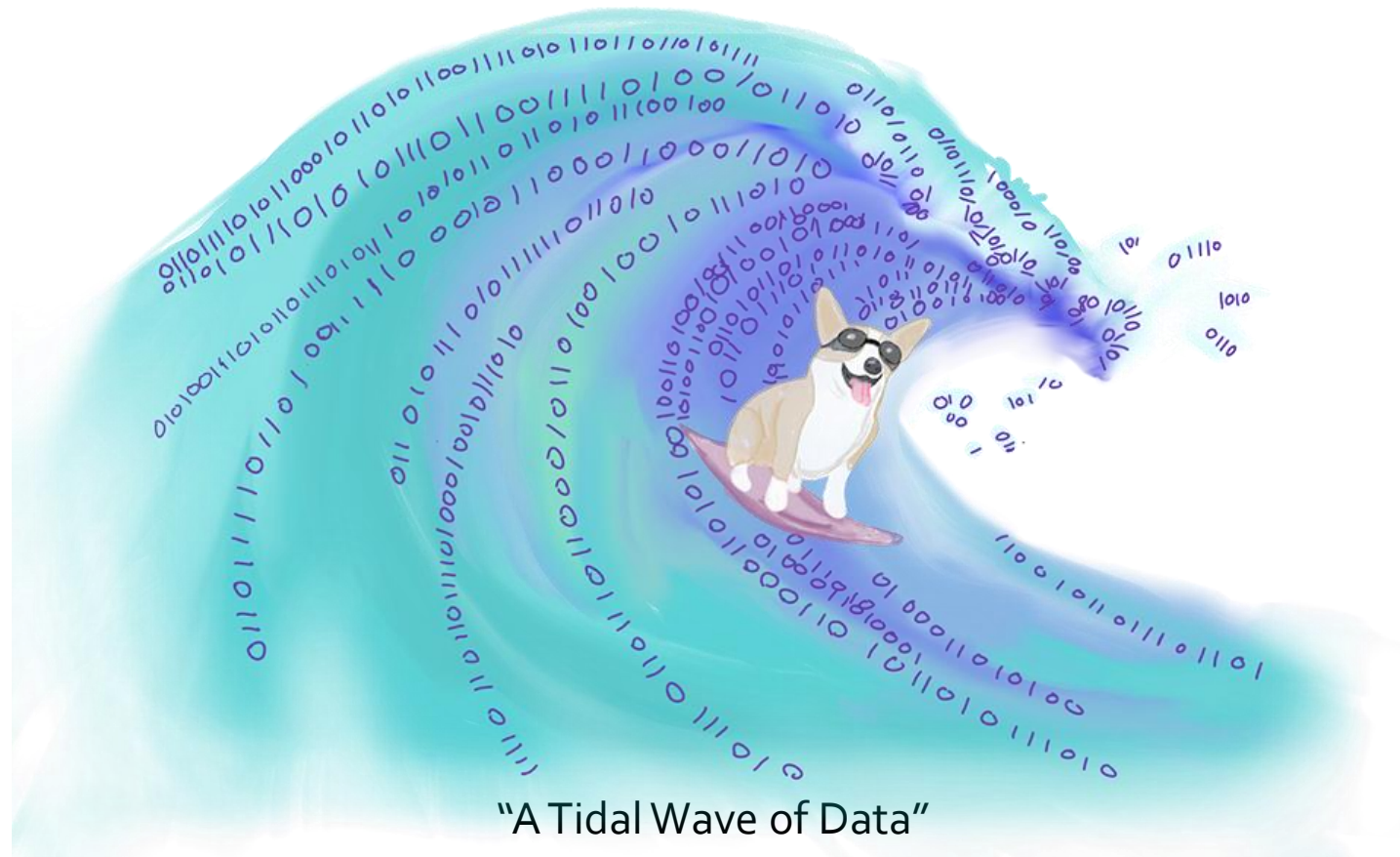
# Motivation ➔ Engagement

# A spectrum

FOR-EACH

WHILE

Games

Websites

Recursion

Iteration

Development

Mobile Apps

IF

Algorithms

Images

**Content**

**Context**

Media
Computation

Assignment

Audio

Data Structures

Animations

Dictionaries

Scientific
Modelling

Scientific
Computing

Lists

Arrays

Integers

Math

Booleans

# Authenticity

- Situated Learning

- "Relevant", "Real-world"

- Media Computation as an "Imagineered Authentic Experience"



*Mark Guzdial and Allison Elliott Tew. 2006. Imagineering inauthentic legitimate peripheral participation: an instructional design approach for motivating computing education. In Proceedings of the second international workshop on Computing education research (ICER '06). New York, NY, USA, 51-58*

# Why *are* we teaching computing?



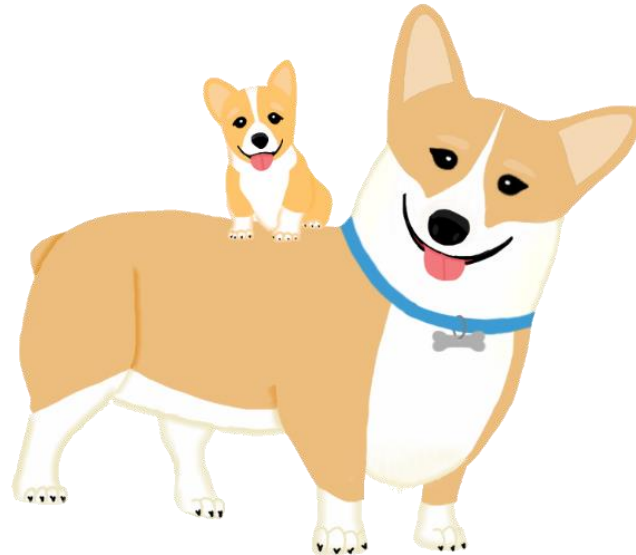"A Tidal Wave of Data"

# State of the Art

- Bart 2014 – Connecting to real-time APIs (RealTimeWeb)

- Hamid 2016 – More generalized framework for real-time APIs (Sinbad)

- Subramanian 2014 – Visualization of data structures with real data (BRIDGES)

- Anderson 2014 – Real world data in CS1

- Sullivan 2013 – Data Science for non-majors

# Problem – We Need Data

- ICPSR – Tightly controlled datasets

- UCI Machine Learning – Only for machine learning

- Census.gov, Kaggle, etc. – Not ready for beginners

# CORGIS

## The Collection Of Really Great, Interesting, Situated Datasets
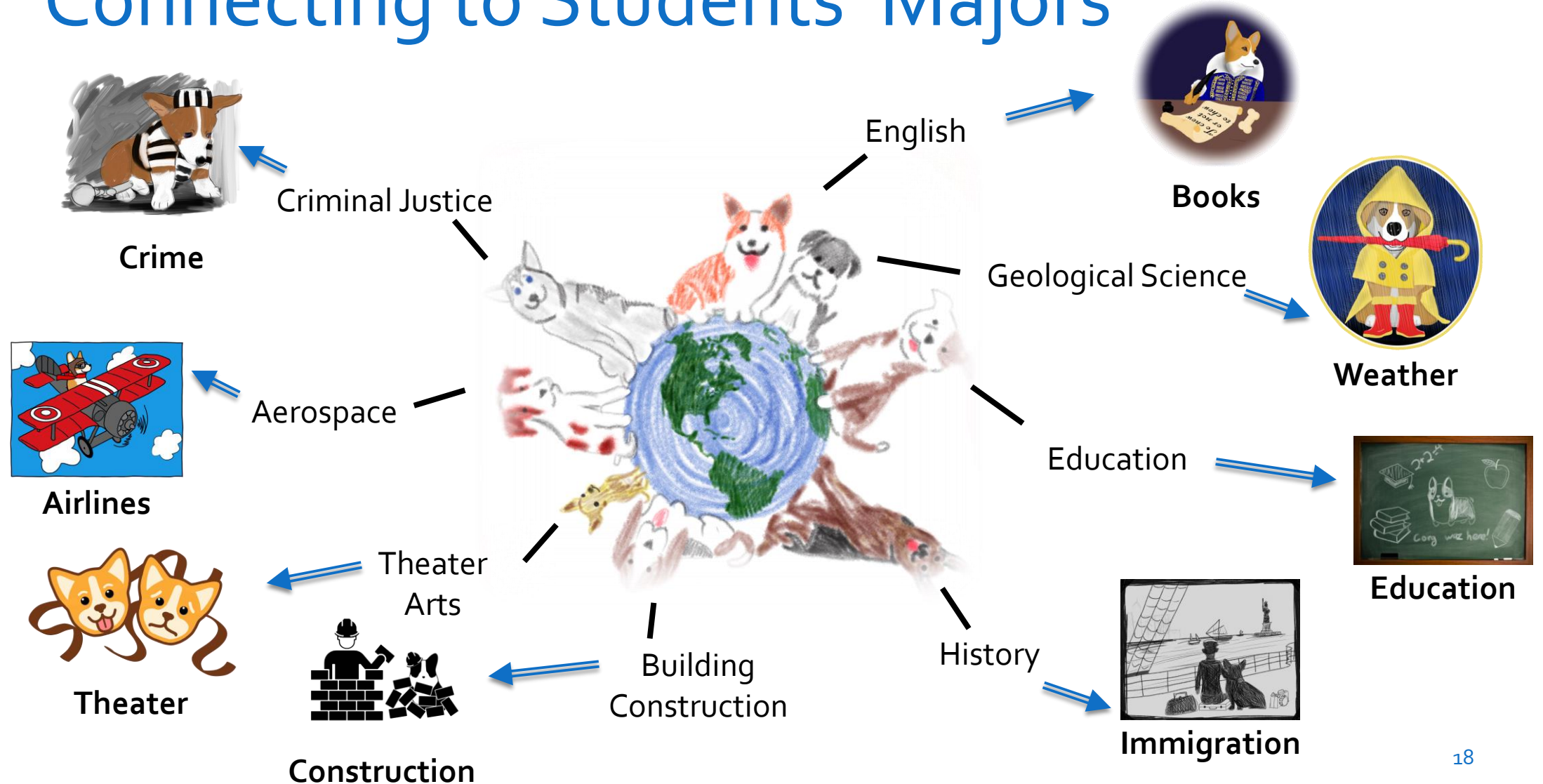
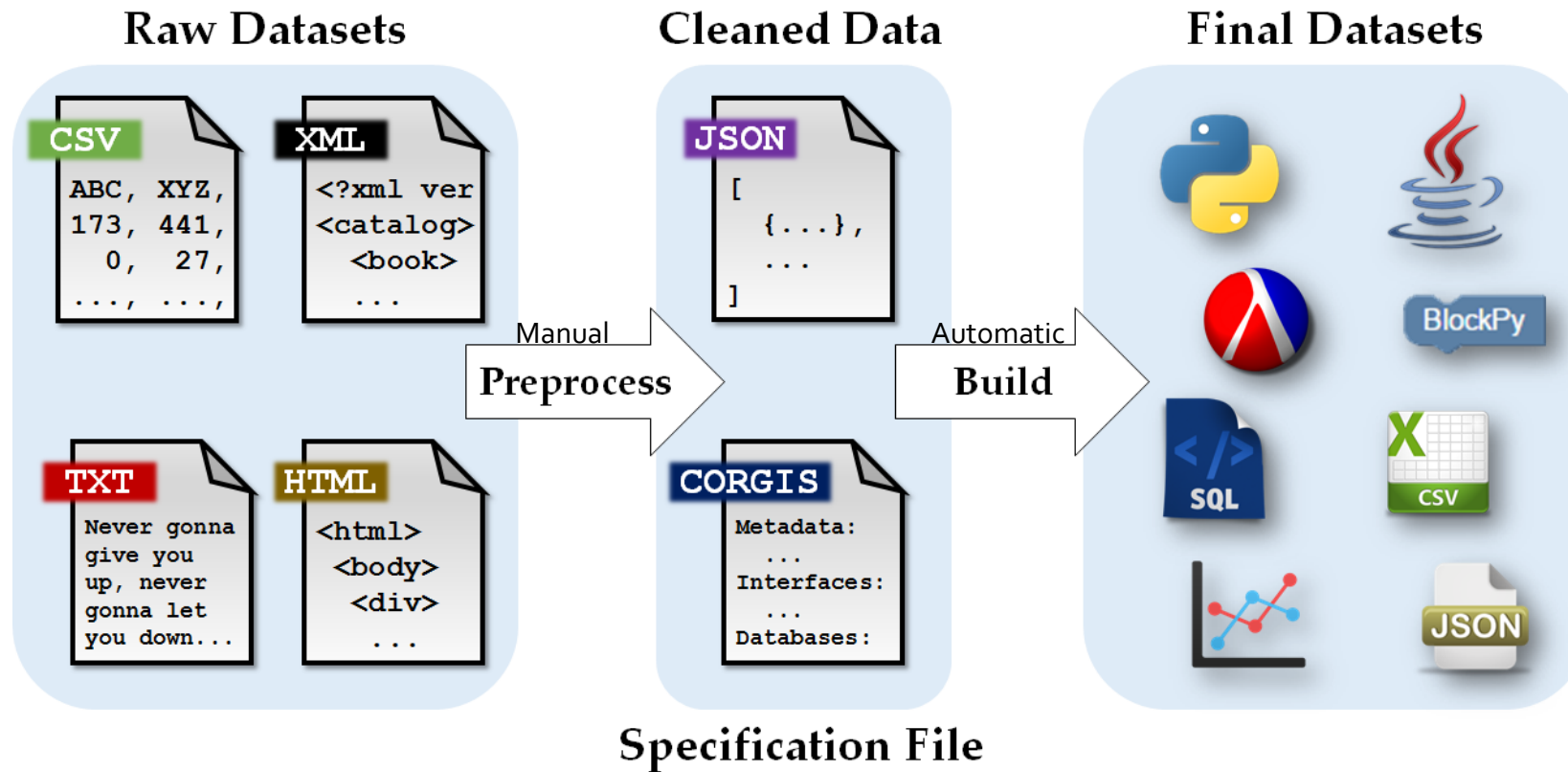# Metrics

42 datasets

267 mB

420,672 rows

9,365,520 values
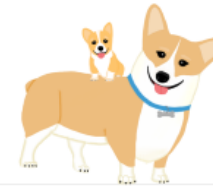
# Datasets

# Connecting to Students' Majors



Criminal Justice — **Crime**

English — **Books**

Geological Science — **Weather**

Aerospace — **Airlines**

Education — **Education**

Theater Arts — **Theater**

Building Construction — **Construction**

History — **Immigration**

# Architecture

# Gallery



## Python Datasets

The **C**ollection of **R**eally **G**reat, **I**nteresting, **S**ituated Datasets

By Austin Cory Bart, Ryan Whitcomb, Jason Riddle, Omar Saleem, Dr. Eli Tilevich, Dr. Clifford A. Shaffer, Dr. Dennis Kafura

**Filter** | Keyword or phrase |

### Aids
Records of AIDS related statistics from several countries.
*aids, death, disease, hiv, orphans, health, countries, world, gender, united nations, un*

### Art Institute Metadata
A data set about the metadata associated with the collection of the Minneapolis Institute of Art.
*art, fine art, institute, artist, style, medium*

### Broadway
This library holds data about Broadway shows, such as tickets sold.
*broadway, musical, theatre, tickets*

### Airlines
Information about flight delays in major aiports since 2003.
*airplane, airports, travel, plane, air, flights, delays, national, united states, transportation*

### Billionaires
Information about over 2000 billionaires from around the world.
*money, rich, wealthy, people, person, billionaire*

### Cancer
Cancer crude rate totals for different ages, races, genders, and geographical areas across the United States.
*cancer, death, states, gender, race, population, crude rate*

# Java, Python, Racket

```java
// Java
import corgis.crime.StateCrimeLibrary;
import corgis.crime.domain.Report;
import java.util.ArrayList;
public class Main {
        public static void main(String[] args) {
                StateCrimeLibrary scl = new StateCrimeLibrary();
                ArrayList<Report> reports = scl.getAll();
        }
}
```

```racket
; Racket
(require crime)
(define reports (crime-get-all))
```

```python
# Python
import crime
crime_reports = crime.get_all()
```

# BlockPy

# Visualizer Demo

# Hypotheses

- Context provides motivation

- Students have some preference for Data Science

- The usefulness of the context connects to engagement outcomes as strongly as the content

# Interventions



- Computational Thinking Course
  - Basic programming
  - Social Impacts
  - Data Science

- 6 semesters taught

- Audience
  - Non-computing majors
  - Freshmen -> Senior
  - Gender balanced

# Motivation × Course Components

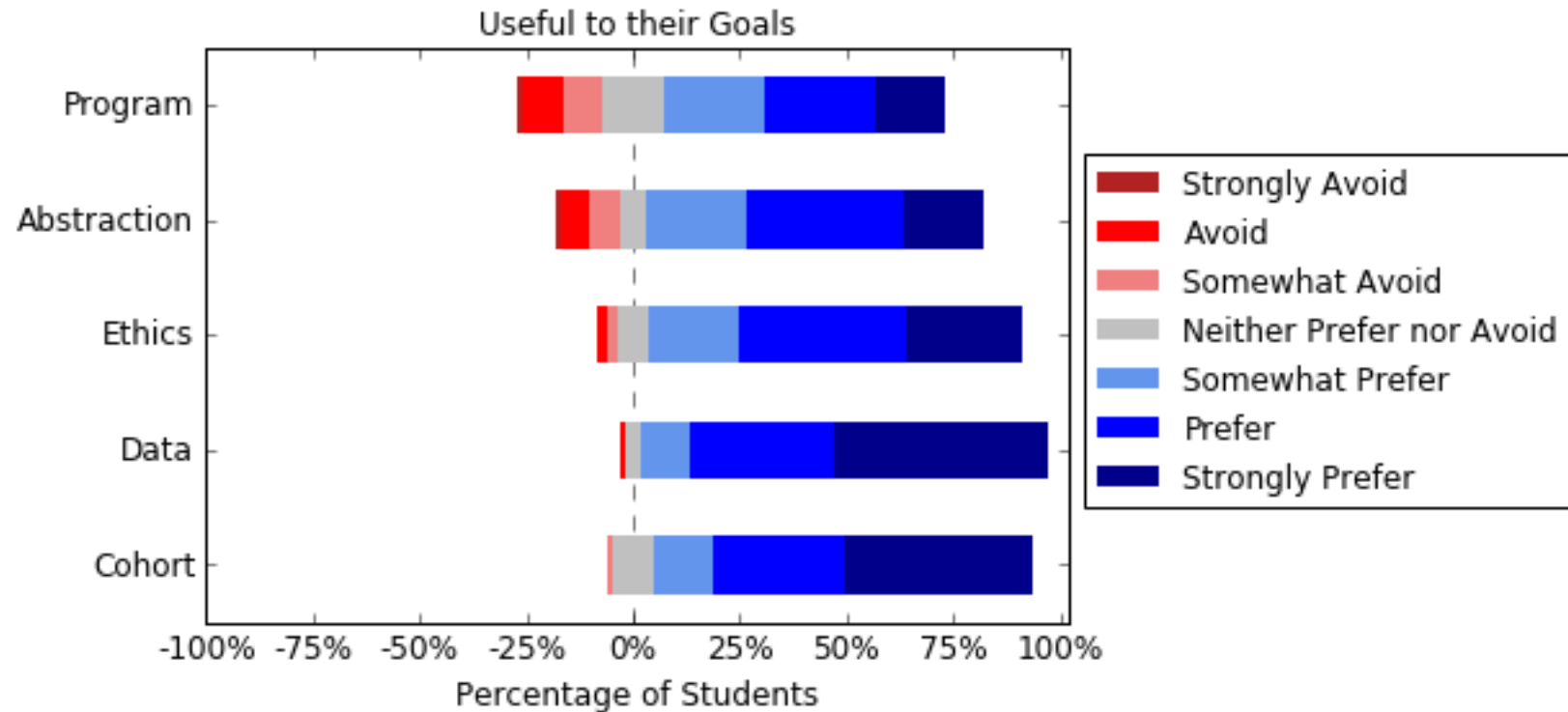| Motivational Components | |
|---|---|
| "I believe that I will have freedom to explore my own interests when I…" | eMpowerment |
| "I believe it will be useful to my long-term career goals to…" | Usefulness |
| "I believe I will be successful in this course when I…" | Success |
| "I believe it will be interesting to…" | Interest |
| "I believe that my instuctors and peers will care about me when I…" | Caring |

| Course Component | |
|---|---|
| "… learn to write computer programs" | Programming Content |
| "… learn to work with abstraction" | Abstraction Content |
| "… learn about the social impacts of computing" | Social Ethics Content |
| "… work with real-world data related to my major" | Data Science Context |
| "… work with my cohort" | Collaboration Facilitation |

| Likert |
|---|
| Strongly Disagree |
| Disagree |
| Somewhat Disagree |
| Neither Agree nor Disagree |
| Somewhat Agree |
| Agree |
| Strongly Agree |

# Context is Useful



Useful to their Goals

N = 85, 62% Female

Students' sense of the usefulness of various course components was highest for the **context**, lowest for the **content**.

# Preference for Contexts

| Preference for Contexts | |
|---|---|
| "Working with data sets related to your major" | Data |
| "Working with pictures, sounds, movies" | Media |
| "Making games and animations" | Games |
| "Making websites" | Web |
| "Making scientific models of real-world phenomenon" | Scientific |
| "Controlling robots or drones" | Robots |
| "Making phone apps" | Mobile |

| Likert |
|---|
| Strongly Avoid |
| Avoid |
| Somewhat Avoid |
| Neither Prefer nor Avoid |
| Somewhat Prefer |
| Prefer |
| Strongly Prefer |

# Preference for Contexts



Student Preference for Introductory Contexts

N = 85, 62% Female

Students' preferred a Data Science context over all others
*No significant difference with Media Computation according to KW test*

# Engagement (Intent to Continue)

| Intent to Continue | |
|---|---|
| "I will try to learn more about computing, either through a course or on my own." | Learn |
| "I will recommend this class to others." | Recommend |
| "I will directly apply what I have learned in my career." | Apply |

| Likert |
|---|
| Strongly Disagree |
| Disagree |
| Somewhat Disagree |
| Neither Agree nor Disagree |
| Somewhat Agree |
| Agree |
| Strongly Agree |

# Engagement (Intent to Continue)



N = 85, 62% Female

Although students would recommend the course, many did not intend to continue learning more computing or applying what they learned.

# Engagement vs. Components

Pearson correlation of "Student's intent to continue learning computing" with students' perception of each course and motivational component

| Fall 2016 | eMpowerment | Usefulness | Success | Interest | Caring |
|---|---|---|---|---|---|
| Abstraction | | | | | |
| Cohort | | | | | |
| Data | | | | | |
| Ethics | | | | | |
| Programming | | .406 | .354 | .341 | |

Not significantly Correlated!

Significant

N = 85, 62% Female

Intent to continue seems to be correlated with the **content**, not the **context**.

# Take-aways

- Data Science seems to be a preferable context for students, across genders.

- Context, and in particular Data Science, can seem to provide motivation in ways that content cannot

- But some engagement outcomes might be more connected to content than context

# Future Work

- More Datasets

- Maintenance

- Connecting motivation to learning outcomes

# Thanks!

Clifford A. Shaffer

Dennis Kafura

Eli Tilevich

Ryan Whitcomb

# Questions?

## https://think.cs.vt.edu/corgis

Artwork by Eleonor Bart

# Trends in Motivation

# Other Components

| Spring 2016 | eMpowerment | Usefulness | Success | Interest | Caring |
|---|---|---|---|---|---|
| Abstraction | .458 | .699 | .614 | .488 | |
| Cohort | | | | | |
| Data | | | | | |
| Ethics | | .485 | .418 | .323 | |
| Programming | .437 | **.823** | .600 | .638 | |

Continue Learning, Applying, and/or Recommend Course
N =36
50% female

# Structure