# Models of support for Data Science
## The perspective of two libraries

CNI Spring Meeting 2024

David Minor, UC San Diego

Jan Brase and Bela Gipp , SUB Göttingen

# Content

1. Introduction
2. UC San Diego, the situation
3. GippLab: Scientific Information Analytics
4. Student and researcher exchange
5. Lessons learned – next steps

# What are we talking about?

Since 2016 the UC San Diego and the State and University Library Goettingen in Germany have been actively cooperating and learned from each other through staff visits and joint projects

MOU (signed in 2016 and renewed in 2019) to compare and contrast research support organizations:

- Comparing different support models
- Relationships between our libraries and respective campuses
- Overall goals and objectives

# Future work (2024 - ?)

At **CNI fall 2023** we presented our longtime Exchanges on RDM infrastructures and services.

For the next years we will additionally focus on:
- Data / research analytics
- Academic exchanges – students
- Further staff exchanges
- Re-examine technical discussions

# UC San Diego, the situation

# UC San Diego Data Science Program - Founded in 2015

Fall 2022: expanded to include three new graduate degrees: in person MS, online MS (joint with Computer Science & Engineering, first fully online graduate degree program offered by UC San Diego)

Fall 2023: ~1000 undergraduate majors, ~6000 students taking courses, ~200 graduate students

*The Data Science Program is primarily focused on education and training students in data science.*

# UC San Diego

## HALICIOĞLU DATA SCIENCE INSTITUTE

**- Founded in 2018**

Brings together faculty members and researchers from various departments across the university to work on data-driven research projects.

Offers various educational programs, including a PhD in Data Science.

Provides resources and support for data science research and education.

*HDSI is focused on advancing the field of data science through research and innovation.*

# Data science in the Library



Hired Data Science Librarian in 2018

- Data & GIS Lab

- Wide range of services

- Bringing student data into Library Digital Collections

- Forthcoming paper in Journal of eScience Librarianship

*Data science support in the Library has to this point been external facing, for students and faculty. We are not yet doing work for internal Library needs.*
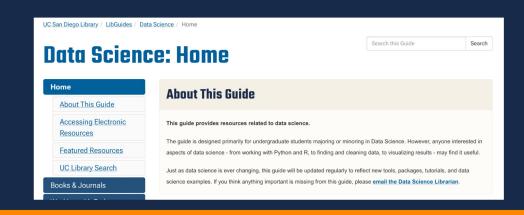
# Data science undergrad support - interactions

- Annual workshops on developing independent (outside of courses) data science projects

- 1-on-1 meetings with students to develop their independent projects, focused on scoping and finding/accessing appropriate data, as well as how to document and present their findings

- Working with capstone student groups to identify data sources for group-based projects

- Working with student groups on their capstone projects - finding and accessing licensed/subscription data; acting as liaison between vendors and student groups when needed

# Data science undergrad support - materials

- Data Science LibGuide with general data science resources (Python tutorials, etc.), and content tied to course curriculum (e.g, assigned readings, direct links to purchased and open access books in library catalog)

- Sharing information about hackathons/datathons, new data resources, Library-hosted events (e.g. UC Love Data Week) and other relevant content in weekly departmental student newsletter

- Quarterly newsletter to faculty about Library updates

- Creating and maintaining a web archive of capstone projects (project websites, GitHub repositories, web-based visualizations, etc.)



UC San Diego Library / LibGuides / Data Science / Home

Search this Guide    Search

## Data Science: Home

**Home**

About This Guide

Accessing Electronic Resources

Featured Resources

UC Library Search

Books & Journals

### About This Guide

**This guide provides resources related to data science.**

The guide is designed primarily for undergraduate students majoring or minoring in Data Science. However, anyone interested in aspects of data science - from working with Python and R, to finding and cleaning data, to visualizing results - may find it useful.

Just as data science is ever changing, this guide will be updated regularly to reflect new tools, packages, tutorials, and data science examples. If you think anything important is missing from this guide, please **email the Data Science Librarian**.

# Data science grad support

- Identifying data suitable for student-organized hackathons/datathons; purchasing requested books for new grad courses; early stages of international student exchanges

- Data Science & Engineering MAS program (not technically HDSI): working directly with students for the annual ingest of group capstone projects

- Computational Social Science grad program (not HDSI, but one of many data science-ish grad programs): working with students to access and analyze data, including leveraging the Library's Data & GIS Lab virtual machines for large data/compute Windows-based analyses

# Next steps: Educational Dataset Service

- Curated datasets for student use in instruction and research:

  - Real research datasets

  - Sanitized and synthetic administrative data (e.g. Student Activity Hub, IR)

  - Datasets commonly used for training in specific disciplines

- Guidance & training for students on working w/associated datasets

  - How to re-use data

  - How to cite sources

  - How to work with metadata

- Support for instructors in developing & assessing meaningful assignments

- Usable in the student Data Science/Machine Learning Platform or elsewhere

# GippLab: Scientific Information Analytics

University of Göttingen and State and University Library (SUB) Göttingen

Bela Gipp

[www.gipplab.org/team](http://www.gipplab.org/team)

# Our Research

Natural Language Processing

Applications

Core Research Areas

Information Retrieval

Natural Language Processing

Blockchain Technology

Foundations

Data Science

Artificial Intelligence

Information Visualization

**Task family setup**

Task family 1 | Task family 2 | ….. | Task family 6

Mixing strategy — proportional / equal

Selection of 18 prepared training tasks

BART

**Training strategies**

Training schemes — sequential / simultaneous / continual multi-task

Pre-fintuned model

finetuning

**Document Classification**

**Dialog Systems**

**Word Sense Disambiguation**

**Machine Translation**

# Natural Language Processing

**Summarization and Generation**

**Question and Answer**

**Sentiment Analysis**

**CoReference Resolution**

---

**Original Text**

...
On **April 29, 2017**, **Bill Gates** partnered with Swiss tennis legend **Roger Federer** in playing the "Match for **Africa**" 4, a noncompetitive tennis match at a sold-out Key Arena in Seattle.
The event was in support of **Roger Federer** Foundation's charity efforts in **Africa**.
...

**Paraphrased using GPT-3**

...
**Bill Gates** teamed up with Swiss tennis player **Roger Federer** to play in the "Match for **Africa** 4" on **April 29, 2017**.
The noncompetitive tennis match at a sold-out Key Arena in Seattle was in support of **Roger Federer** Foundation's charity efforts in **Africa**.
...

---

⚡ **Hosted inference API** ⓘ

Text2Text Generation          Examples ⌄

question: which description describes the word " java " best in the following context? descriptions: [ " A drink consisting of an infusion of ground coffee beans " , " a platform-independent programming lanugage " , or " an island in Indonesia to the south of Borneo " ] context: I like to drink ' java ' in the morning .

Compute          ⌘+Enter          1.7

Computation time on cpu: 1.596 s

a drink consisting of an infusion of ground coffee beans

</> JSON Output          ⤢ Maximize

🖲 Space using jpwahle/t5-word-sense-disambigua…

📄 linpearjun/T5-LARGE

---

TF 1 — Task 1
TF 2 — Task 3
TF 6 — Task 2
TF 1 — Task 3
TF 1 — Task 1

(a) Sequential learning.

TF 1 | TF 2 | TF 3
TF4 | TF 5 | TF 6

(b) Simultaneous learning.

TF 1 — Task 1
TF 1 — Task 2 | TF 1 — Task 1
TF 1 — Task 3 | TF 1 — Task 1 | TF 1 — Task 2
TF 2 — Task 1 | TF 1 — Task 1 | TF 1 — Task 2 | TF 1 — Task 3
TF 2 — Task 2 | TF 1 — Task 1 | TF 1 — Task 2 | TF 1 — Task 3 | TF 2 — Task 1

(c) Continual multi-task learning.

$d_s$= Albert Einstein

$r_1$= Country of Citzenship

$r_2$=Educated at

$d_{t1}$= German Empire

$d_{t2}$= ETH Zurich

$r_1$

$r_2$

*Other persons who have the same citizenship*

Otto von Bismarck

Wilhelm I

Max von Baden

...

*Other persons who were educated at the same instituion*

Wilhelm Conrad Röntgen

Charles-Edouard

Guillaume

Otto Stern

...

newer documents

older documents

[1]
[2]

Doc B

[3]

inspired by

[1]
[2]

Doc A

[3]

cites

[1]

[2]

[3]

Note: Doc A receives no citation from Doc B

# Literature Recommendation

.43
Performance Characteristics of the Abbot Architect SARS-CoV-2 IgG Assay and Seraprevalence in Boise, Idaho

.65
Hypertension and Renin-Angiotensin-Aldosterone System Inhibitors in Patients with Covid-19

.85
Observational Study of Hydroxychloroquine in Hospitalized Patients with Covid-19

.55
Impact of blood analysis and immune function on the prognosis of patients with COVID-19

.35

.17
Observational Study of Hydroxychloroquine in Hospitalized Patients with Covid-19

A Trial of Lopinavir–Ritonavir in Adults Hospitalized with Severe Covid-19

Cao et al.

New England Journal of Medicine

7 May, 2020

inspect

GLOBAL SENSITIVITY THRESHOLD

SIMILARITY WEIGHT

TEXT

CITATIONS

FIGURES

CUSTOM KEYWORDS

save weights

SIMILARITY WEIGHT SELECTION

○ custom citations
◉ full
○ zero

18

# Plagiarism Detection

1. User submits the unique fingerprint (SHA-256) of his file to OriginStamp service

2. Once a day, OriginStamp aggregates the submitted hashes and embeds them into a transaction that is broadcasted over the network

3. Anyone can verify the tamper-proof timestamp using the distributed blockchain once the transaction is confirmed by the public, distributed blockchain

SUBMIT

HASH

STORE

VERIFY

9 d 9 1 4 b e 6 9 7 d

# IP Protection

Timestamp your documents

Location-based Exploration

Prove the existence of documents

# News Analysis



Objective

Frame 1

Frame 2

# Math Information Retrieval

Applications

Applications

**Core Research Areas**

Information Retrieval

Natural Language Processing

Blockchain Technology

**Foundations**

Data Science

Artificial Intelligence

Information Visualization
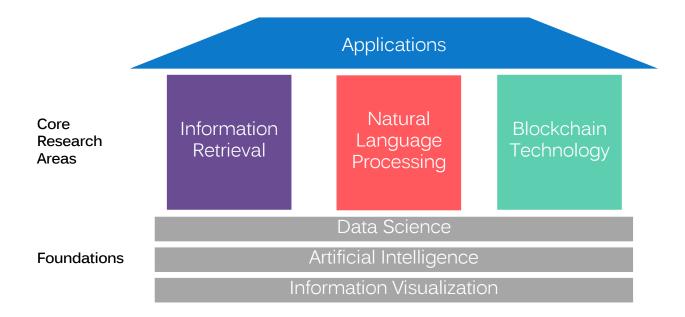
# Researcher and Student exchange

- Academic exchanges – students
    - UCSD students mandatory have to work on a project for their data science master - SUB and University of Goettingen have many projects
    - Uni Goettingen students are encouraged to spend a semester abroad and work on projects – Extrensiv motivtation for excellent students.

- Further staff exchanges
    - Walk in each others shoes, looking over the shoulder in understanding daily work
    - Networking

# Funding Programms for Exchange

There are several scholarships by the DAAD for:

- Undergraduates: [LINK](#)

- Graduates, PhDs, Postdocs from the US
  - up to 6 months: [LINK](#)
  - up to 12 months: [LINK](#)
- For jointly supervised PhD candidates: [LINK](#)

- For academics & scientists (up to 3 months): [LINK](#)



SCAN ME

# Lessons learned- Next steps

UCSD library and SUB Göttingen have benefitted from learning from each other through staff visits.

Additionally UCSD and University of Göttingen will cooperate in data science research and teaching, this time fueled by researchers and student exchanges.

We aim at having the first exchanges in 2024-2025

# Thank you

Jan Brase - [brase@sub.uni-goettingen.de](mailto:brase@sub.uni-goettingen.de)
David Minor – [dminor@ucsd.edu](mailto:dminor@ucsd.edu)
Bela Gipp- [gipp@uni-goettingen.de](mailto:gipp@uni-goettingen.de)
[www.gipplab.org](http://www.gipplab.org)