# ENGAGING WITH THE RESEARCH COMMUNITY

AI and Data Transformation

John Little ⓘD

Duke University Libraries

Center for Data & Visualization Sciences

2024-10-27

# PATRON-DRIVEN INTERACTIONS

Computational workflows, data transformations, and Analysis

# AI AND NEW CHALLENGES

- Introduction of Generative LLMs (e.g., ChatGPT)

- New challenges for patrons:

  1. How to ask questions of a generative AI

  2. How to frame questions to reflect data goals

- Translation

- Synethesis

# THE CONFIDENCE V COMPETENCE PARADOX

- LLMs give confident responses

- Responses are predictions, not necessarily correct answers

- Incorrect predictions = "hallucinations"

- Verification is crucial

- Paradox: More knowledge leads to better evaluation of AI responses

# USE CASE - CODE GENERATION

- Data transformation

- Data analysis

- Iteration

- Big Data

- AI assistance / AI-paired coding

# GOAL

Create scatter plots, one for each home world

# CASE STUDY - STAR WARS DATASET

| Homeworld | Heights | Masses | Characters |
|---|---|---|---|
| Tatooine | 172, 188, 178 | 77, 84, 120 | Luke Skywalker, Anakin Skywalker, Owen Lars |
| Alderaan | 150, 191 | 49, 85 | Leia Organa, Bail Prestor Organa |
| Naboo | 165, 196, 170 | 45, 66, 75 | Padmé Amidala, |

| Homeworld | Heights | Masses | Characters |
|---|---|---|---|
| | | | Isadore, Palpatine |
| Coruscant | 66, 188 | 17, 84 | Yoda, Mace Windu |

# EXAMPLE

**Star Wars Characters from Alderaan**
Mass vs Height



**Star Wars Characters from Tatooine**
Mass vs Height



**Star Wars Characters from Naboo**
Mass vs Height

# CHALLENGES IN AI ASSISTANCE

- AI *can* handle well some basic visualization and coding

- Struggles with complex data shaping and iteration

- This problem is easier when the user has knowledge in:

  - Coding concepts

  - Data shaping

  - Visualization

  - Iteration for large datasets

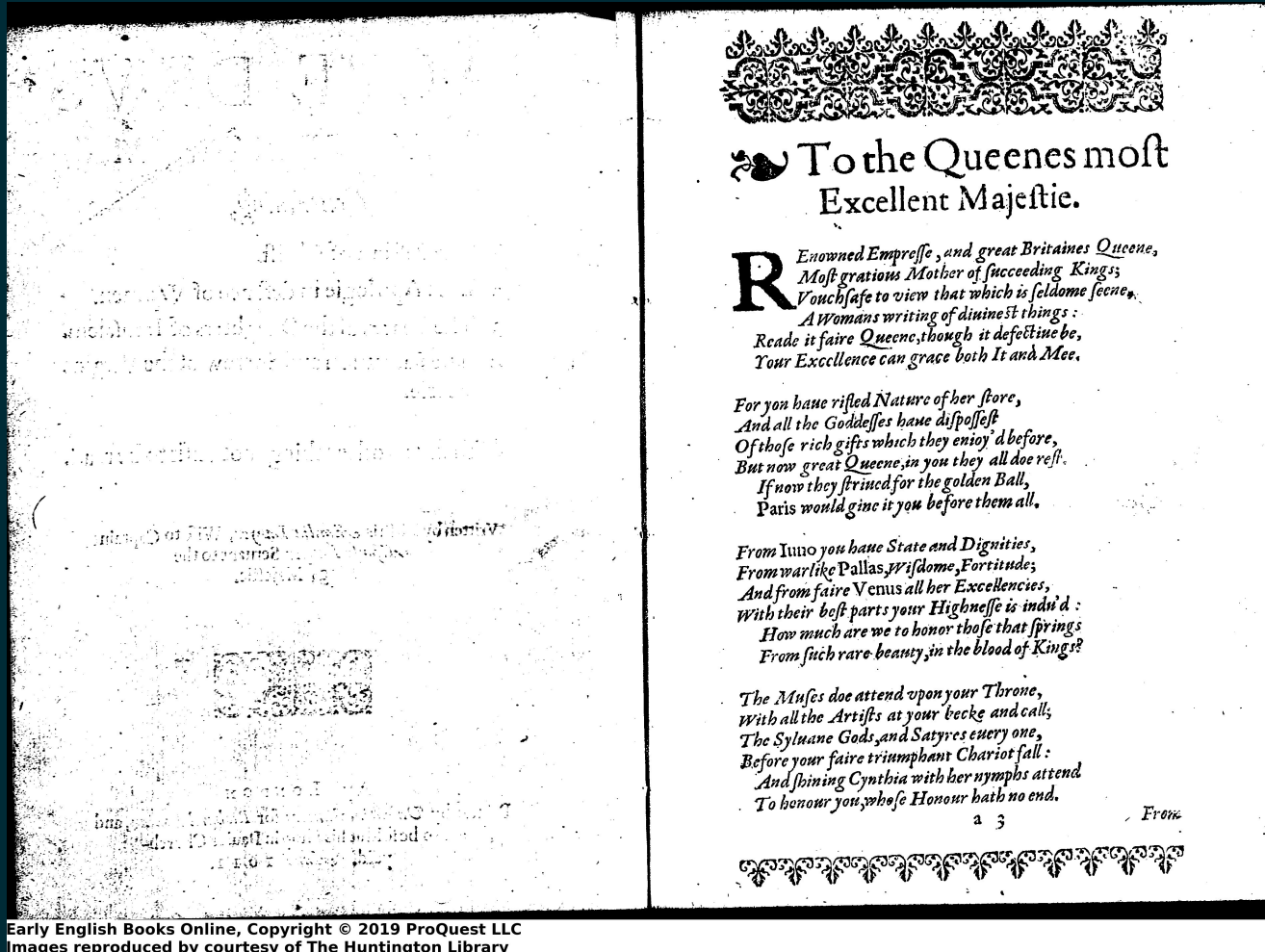# WHEN IT GOES WRONG

# WORD PROBLEMS

Prompt: Inconsistent AI responses for "How long does it take to walk 10,000 steps on a treadmill at 1.2 MPH?"

- Lesson 1: Importance of cross-verification

- Lesson 2: Prediction is not the same as mathmatical truth

# EEBO

## No ground truth



**To the Queenes most** Excellent Majestie.

Renowned Empresse, and great Britaines Queene,
Most gratious Mother of succeeding Kings;
Vouchsafe to view that which is seldome seene,
A Womans writing of diuinest things:
Reade it faire Queene, though it defectiue be,
Your Excellence can grace both It and Mee.

For you haue rifled Nature of her store,
And all the Goddesses haue dispossest
Of those rich gifts which they enioy'd before,
But now great Queene, in you they all doe rest.
If now they striued for the golden Ball,
Paris would giue it you before them all.

From Iuno you haue State and Dignities,
From warlike Pallas, Wisdome, Fortitude;
And from faire Venus all her Excellencies,
With their best parts your Highnesse is indu'd:
How much are we to honor those that springs
From such rare beauty, in the blood of Kings?

The Muses doe attend vpon your Throne,
With all the Artists at your becke and call;
The Syluane Gods, and Satyres euery one,
Before your faire triumphant Chariot fall:
And shining Cynthia with her nymphs attend
To honour you, whose Honour hath no end.

a 3                    From

# CODE

## TRANSLATION DONE POORLY

- Due to insufficient background and/or prompting

## AI-PAIRED CODE GENERATION

- Some clear winners and losers in the big names. aka each LLM has it's own evolving strengths, weaknesses, and tendencies.

These problem highlights the Competence v Confidence Paradox but are easily **verifiable**

# WHEN IT GOES RIGHT

and how *right* does it go?

# SYNETHTIC QUESTIONS

Prompt:   Compare student body and faculty diversity at Duke University with UNCG. Compare today with 1985.

- Lesson 1: Different LLMs give different amounts of evidence for verification

- Lesson 2: Differing amounts of ground truth will affect the prediction

# CODE TRANSLATION

I have Python code, give it to me in R

# VARIATIONS IN CODE TRANSLATIONS

- R to Python

- Python to R

- SQL from natural language

- javascript

- HTML

# NATURAL LANGUAGE

How can I use the phrase "*Sticky Wicket*" in German?

- Translate *Sticky Wicket* to German

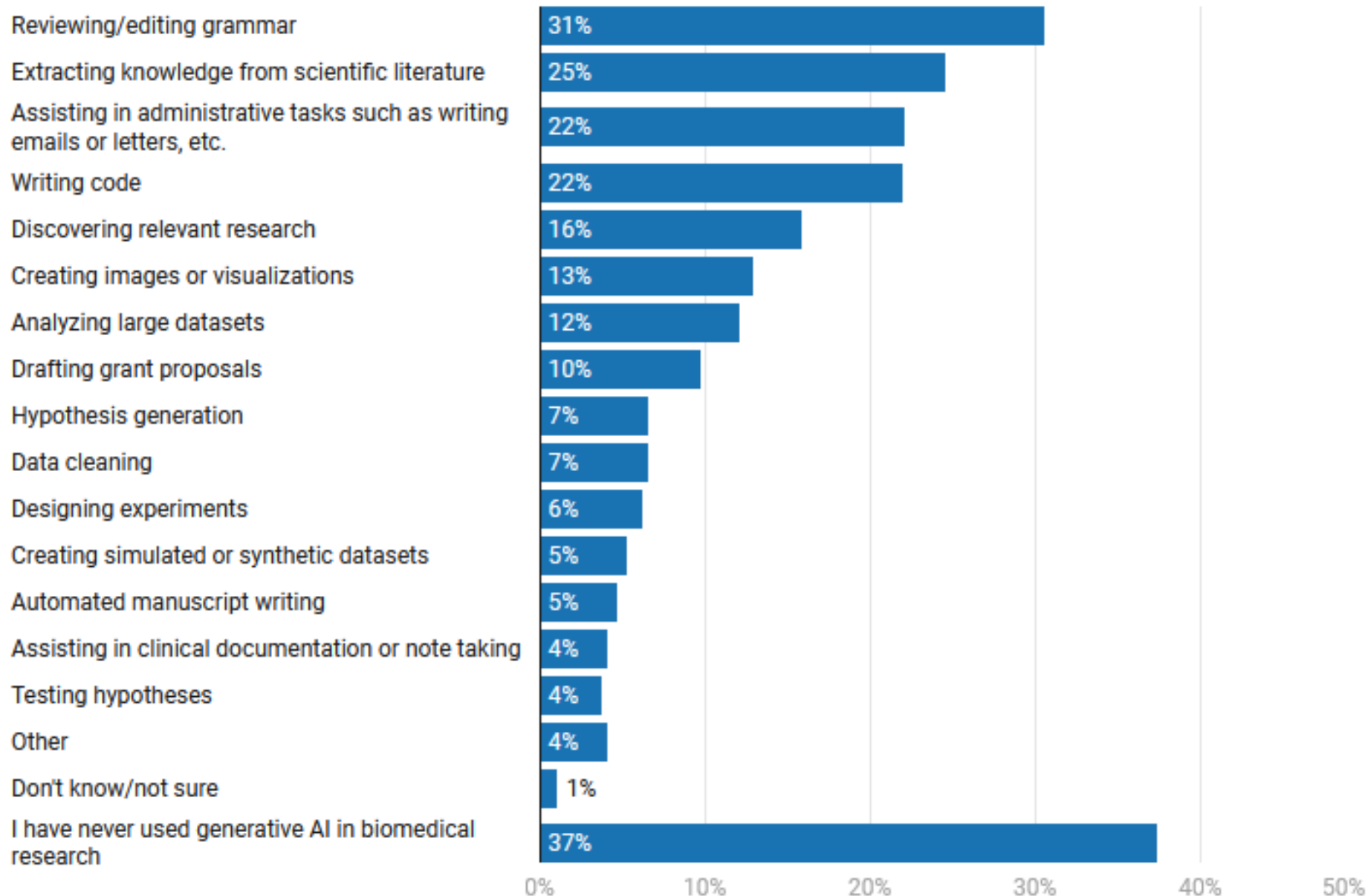- But how to verify (same as code problem)

# VALUE IN REPRODUCIBILITY

- Coding

  - Do everything with code

  - Including report generation

- No Code

  - Getting better all the time

Increasingly we are seeing computation environments with build-in AI-pairing

**Figure 9. "In which of the following ways have you used generative AI in your biomedical research?"**

| Category | Percentage |
|---|---|
| Reviewing/editing grammar | 31% |
| Extracting knowledge from scientific literature | 25% |
| Assisting in administrative tasks such as writing emails or letters, etc. | 22% |
| Writing code | 22% |
| Discovering relevant research | 16% |
| Creating images or visualizations | 13% |
| Analyzing large datasets | 12% |
| Drafting grant proposals | 10% |
| Hypothesis generation | 7% |
| Data cleaning | 7% |
| Designing experiments | 6% |
| Creating simulated or synthetic datasets | 5% |
| Automated manuscript writing | 5% |
| Assisting in clinical documentation or note taking | 4% |
| Testing hypotheses | 4% |
| Other | 4% |
| Don't know/not sure | 1% |
| I have never used generative AI in biomedical research | 37% |

John R. Little • Center for Data & Visualization Sciences • CC BY 4.0

# SOLUTIONS

and best practices

# PROBLEMS AND SOLUTIONS

- GIGO (Garbage In, Garbage Out) still applies

- Prompt engineering is a crucial skill

- AI excels in translation tasks

- Good for synthetic questions with possible validation

- Less reliable for tasks without established ground truth

# BEST PRACTICES

Using Broad-base LLMs:

- ChatGPT

- Microsoft Copilot

- Claude.ai

- Gemini.google.com

- GitHub Copilot (for AI-paired coding)

# PROMPT ENGINEERING

- Identify role

- Identify audience

- Identify voice

- Identify goals and problem

- Use multiple steps

- Verify

# CONCLUSION

Embracing AI in data analysis

- AI is a powerful tool, but requires careful use

- The library offers crucial guidance

- Continuous learning and adaptation are essential

# QUESTIONS

1. How do you see these tools or techniques impacting research and research investment?
2. Do you have data transrormation, reshaping, or analysis tasks that could benefit from AI assistance?
3. In what ways do you think we can improve training and assistance for next generation LLMs?
4. What are some of the biggest challenges you see in the future of AI-paired coding?