

# NLM Scrubber to de-identify clinical text data and facilitate data sharing

---



**Lightning Talk – MIRL 2023**

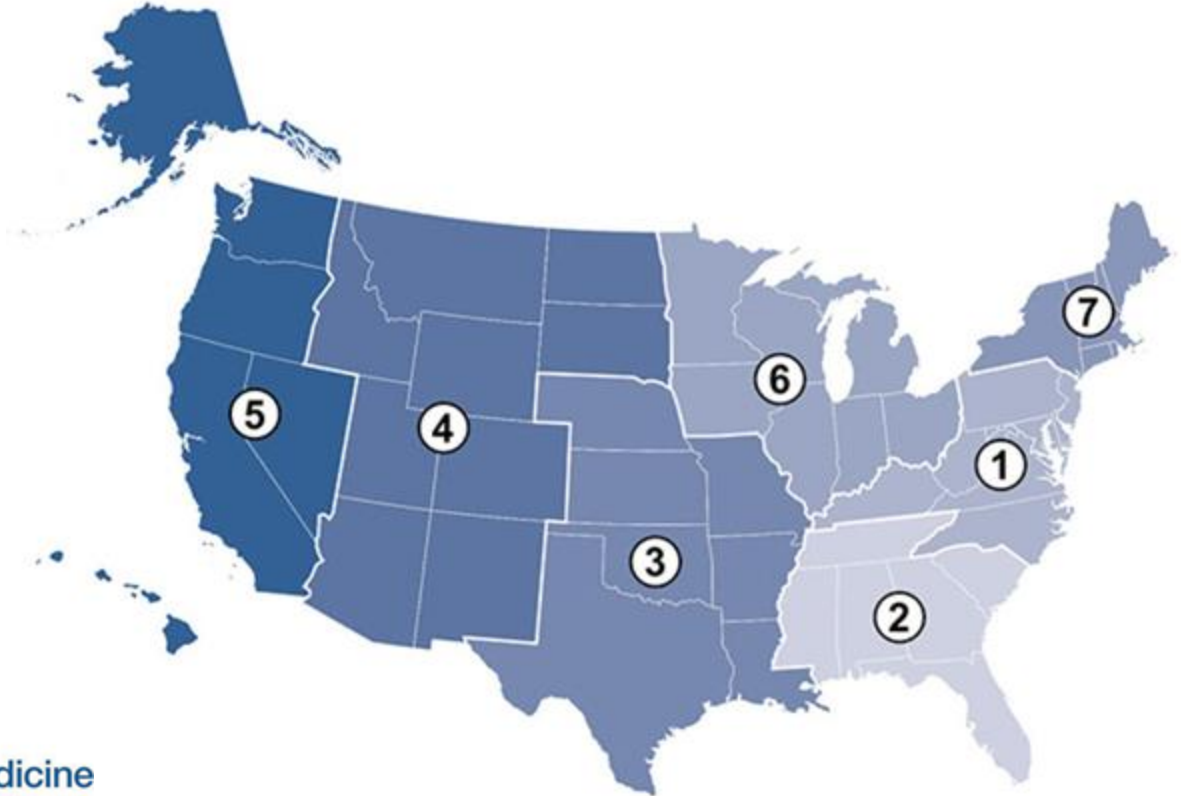
Katie Pierce Farrier, Data Science Strategist

Christine Nieman Hislop, Data Education Librarian

# Network of the National Library of Medicine

---

- Education
- Outreach
- Funding



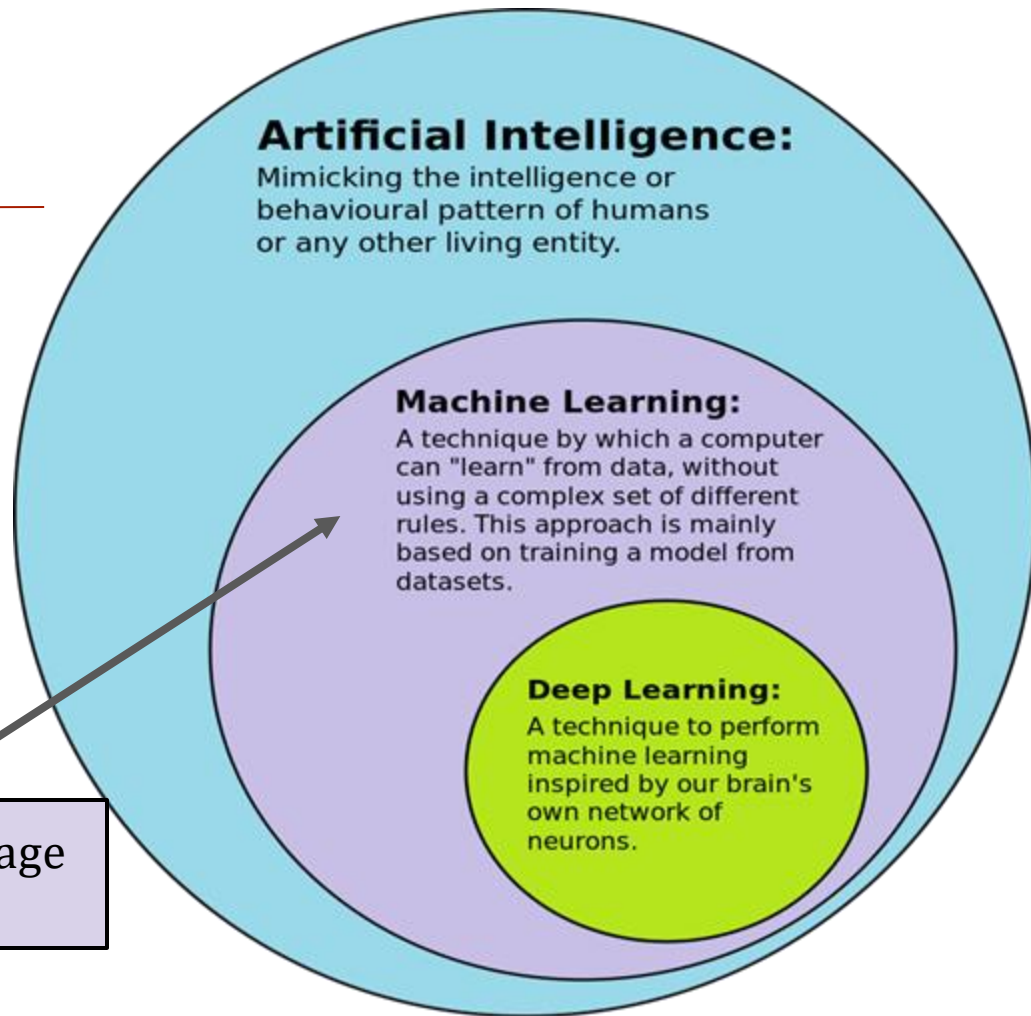
# Defining AI

---

*No megalomaniac robots  
involved....for now*

[Photo from Wikimedia Commons](#)

Natural Language  
Processing



# So what is NLP?

---

- Languages follow predictable **rules**
- Machine Learning follows **statistics**
- Allows computers to **understand human languages**



# Examples of NLPs

- Apple Siri
- Amazon Alexa
- ChatGPT
- Customer service chatbots

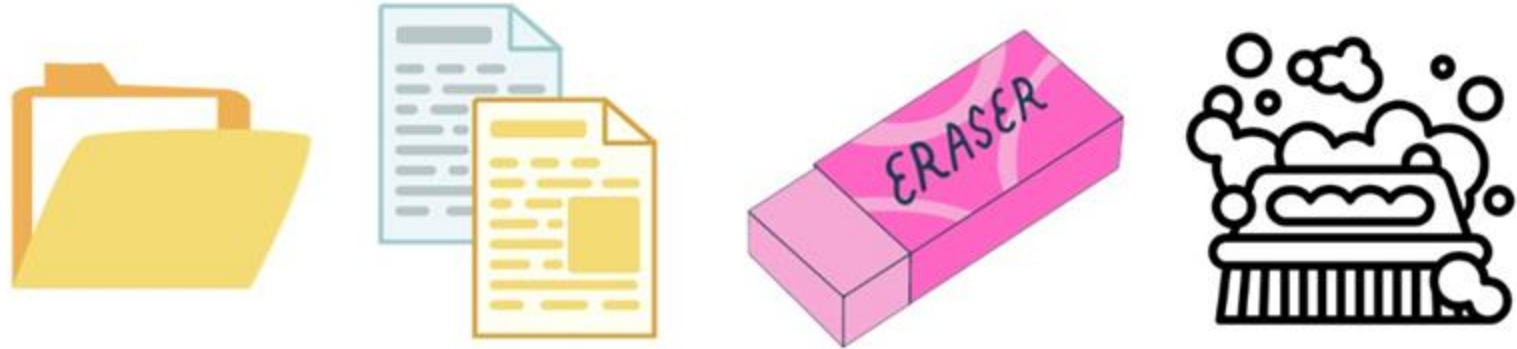


National Library of Medicine  
Network of the National Library of Medicine



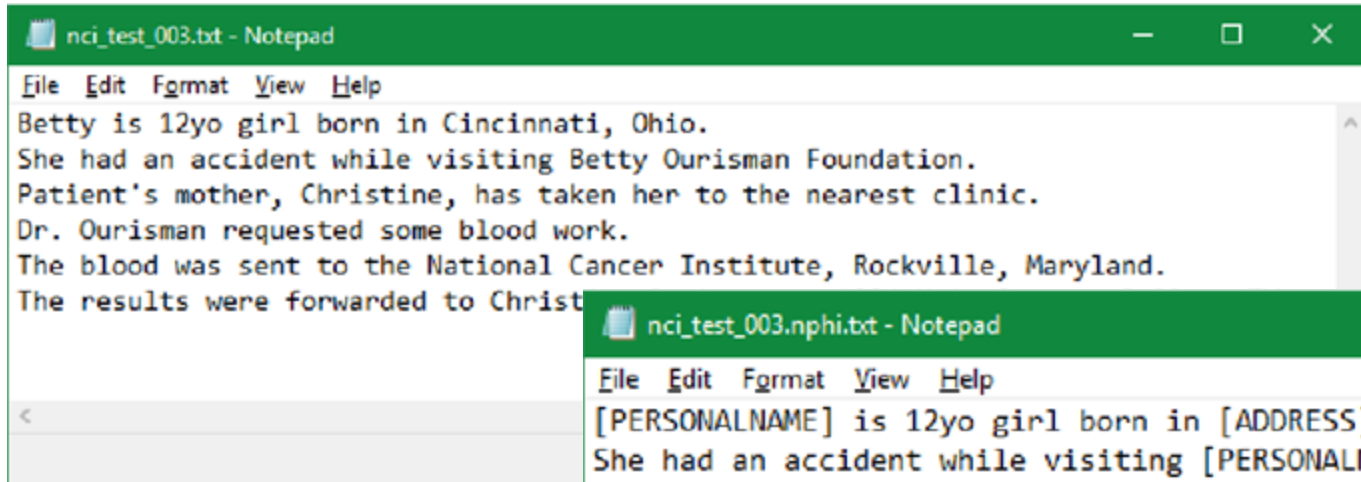
# NLM Scrubber uses NLP to find and redact PHI ....

---

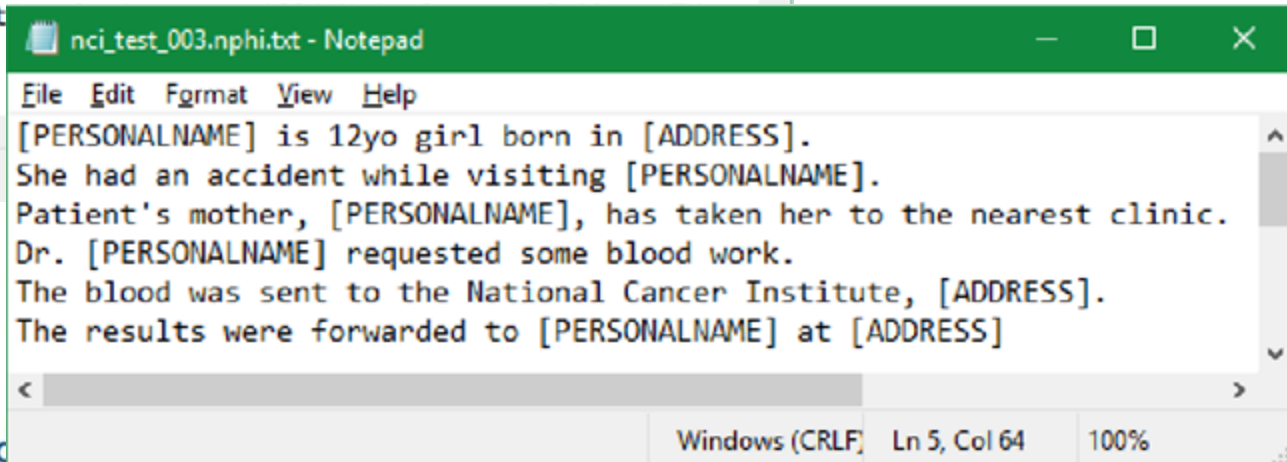


# ... to protect patient privacy and...

---



```
nci_test_003.txt - Notepad
File Edit Format View Help
Betty is 12yo girl born in Cincinnati, Ohio.
She had an accident while visiting Betty Ourisman Foundation.
Patient's mother, Christine, has taken her to the nearest clinic.
Dr. Ourisman requested some blood work.
The blood was sent to the National Cancer Institute, Rockville, Maryland.
The results were forwarded to Christ
```



```
nci_test_003.nphi.txt - Notepad
File Edit Format View Help
[PERSONALNAME] is 12yo girl born in [ADDRESS].
She had an accident while visiting [PERSONALNAME].
Patient's mother, [PERSONALNAME], has taken her to the nearest clinic.
Dr. [PERSONALNAME] requested some blood work.
The blood was sent to the National Cancer Institute, [ADDRESS].
The results were forwarded to [PERSONALNAME] at [ADDRESS]
Windows (CRLF) Ln 5, Col 64 100%
```



# .... facilitate sharing of clinical de-identified data!

---





# De-identified Data = Shareable Data

---



**F**indable



**A**ccessible



**I**nteroperable



**R**eusable

Don't forget **Reproducible!**

- 68% of experiments unable to obtain data
- 39 of 53 papers cited NIH funding (Errington et al., 2021)





# Thank you!

---



Christine Nieman Hislop  
[cnieman@hshsl.umaryland.edu](mailto:cnieman@hshsl.umaryland.edu)  
NNLM Region 1



Katie Pierce Farrier  
[Katie.pierce-farrier@unthsc.edu](mailto:Katie.pierce-farrier@unthsc.edu)  
NNLM Region 3



# Resources

---

[NLM Scrubber](#) and its [User Manual](#)

[IBM: What is Natural Language Processing?](#)

[Go FAIR Data Principles](#)

[Challenges for assessing replicability in preclinical cancer biology](#)

## Articles About NLM Scrubber

[Evaluation of Automated Public De-Identification Tools on a Corpus of Radiology Reports](#)

[An atomic approach to the design and implementation of a research data warehouse](#)

[A Comparative Analysis of Speed and Accuracy for Three Off-the-Shelf De-Identification Tools](#)

