



# With great power comes great responsibility:

## Machine learning in clinical research

*E. Hope Weissler, MD, MHS & Erich Huang, MD, PhD*

# Disclosures

**Weissler:** None

**Huang:**

Founder - kelaHealth, Clinetic, Stratus Medicine  
Chief Science and Innovation Officer at Onduo (a Verily company)  
Advisor – Optimizely

# Contents



## Past

*What is holding clinical  
research back from  
reaching full potential?*

# Contents



## Past

*What is holding clinical research back from reaching full potential?*

## Present

*How might machine learning address those issues?*

*What are the barriers to ML implementation in CR?*



# Contents



## Past

*What is holding clinical research back from reaching full potential?*

## Present

*How might machine learning address those issues?*

*What are the barriers to ML implementation in CR?*

## Future

*How can we overcome these barriers?*

# Contents

## DCRI THINK TANKS

FROM INSIGHT TO ACTION

***Leveraging Artificial Intelligence and Machine Learning Methods and Approaches to Transform Clinical Trial Design, Planning, and Execution***

**~50 attendees including:**

- clinical researchers
- machine learning experts
- biopharmaceutical industry
- technology companies
- patient advocacy groups
- FDA

Weissler et al. *Trials* (2021) 22:537  
<https://doi.org/10.1186/s13063-021-05489-x>

Trials

COMMENTARY

Open Access

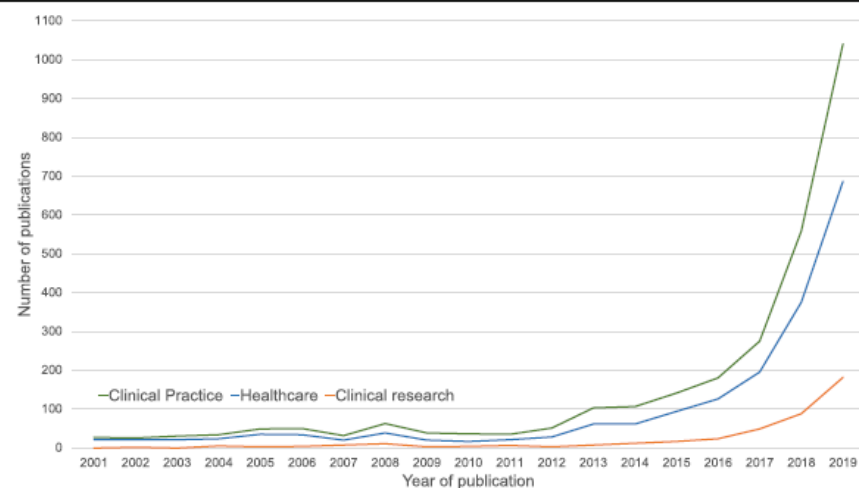
### The role of machine learning in clinical research: transforming the future of evidence generation



E. Hope Weissler<sup>1\*</sup>, Tristan Naumann<sup>2</sup>, Tomas Andersson<sup>3</sup>, Rajesh Ranganath<sup>4</sup>, Olivier Elemento<sup>5</sup>, Yuan Luo<sup>6</sup>, Daniel F. Freitag<sup>7</sup>, James Benoit<sup>8</sup>, Michael C. Hughes<sup>9</sup>, Faisal Khan<sup>3</sup>, Paul Slater<sup>10</sup>, Khader Shameer<sup>3</sup>, Matthew Roe<sup>11</sup>, Emmette Hutchison<sup>3</sup>, Scott H. Kollins<sup>1</sup>, Uli Broedl<sup>12</sup>, Zhaoling Meng<sup>13</sup>, Jennifer L. Wong<sup>14</sup>, Lesley Curtis<sup>1</sup>, Erich Huang<sup>1,15</sup> and Marzyeh Ghassemi<sup>16,17,18,19</sup>

Ref 1.

# Present: What can we do with ML in CR?




**Fig. 1** The number of clinical practice-related publications was determined by searching “(“machine learning” or “artificial intelligence”) and (“healthcare”).” The number of healthcare-related publications was determined by searching “(“machine learning” or “artificial intelligence”) and (“healthcare”).”, and the number of clinical research-related publications was determined by searching “(“machine learning” or “artificial intelligence”) and (“clinical research”).”

Weissler et al. *Trials* (2021) 22:537  
<https://doi.org/10.1186/s13063-021-05489-x>

## BRIEF COMMUNICATION

OPEN

# The National Institutes of Health funding for clinical research applying machine learning techniques in 2017

Amarnath R. Annapureddy <sup>1,2</sup>, Suveen Angraal<sup>1,3</sup>, Cesar Caraballo <sup>1</sup>, Alyssa Grimshaw <sup>4</sup>, Chenxi Huang<sup>1</sup>, Bobak J. Mortazavi<sup>5</sup> and Harlan M. Krumholz <sup>1,2,6\*</sup>

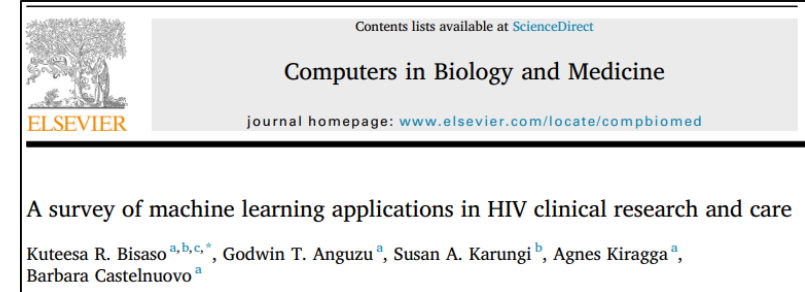
# Contents

Prior  
reviews of  
the topic

Circulation: Cardiovascular Quality and Outcomes

NOVEL STATISTICAL METHODS

**Recommendations for Reporting Machine Learning Analyses in Clinical Research**



Ref. 3-5

**Pre-clinical  
work**

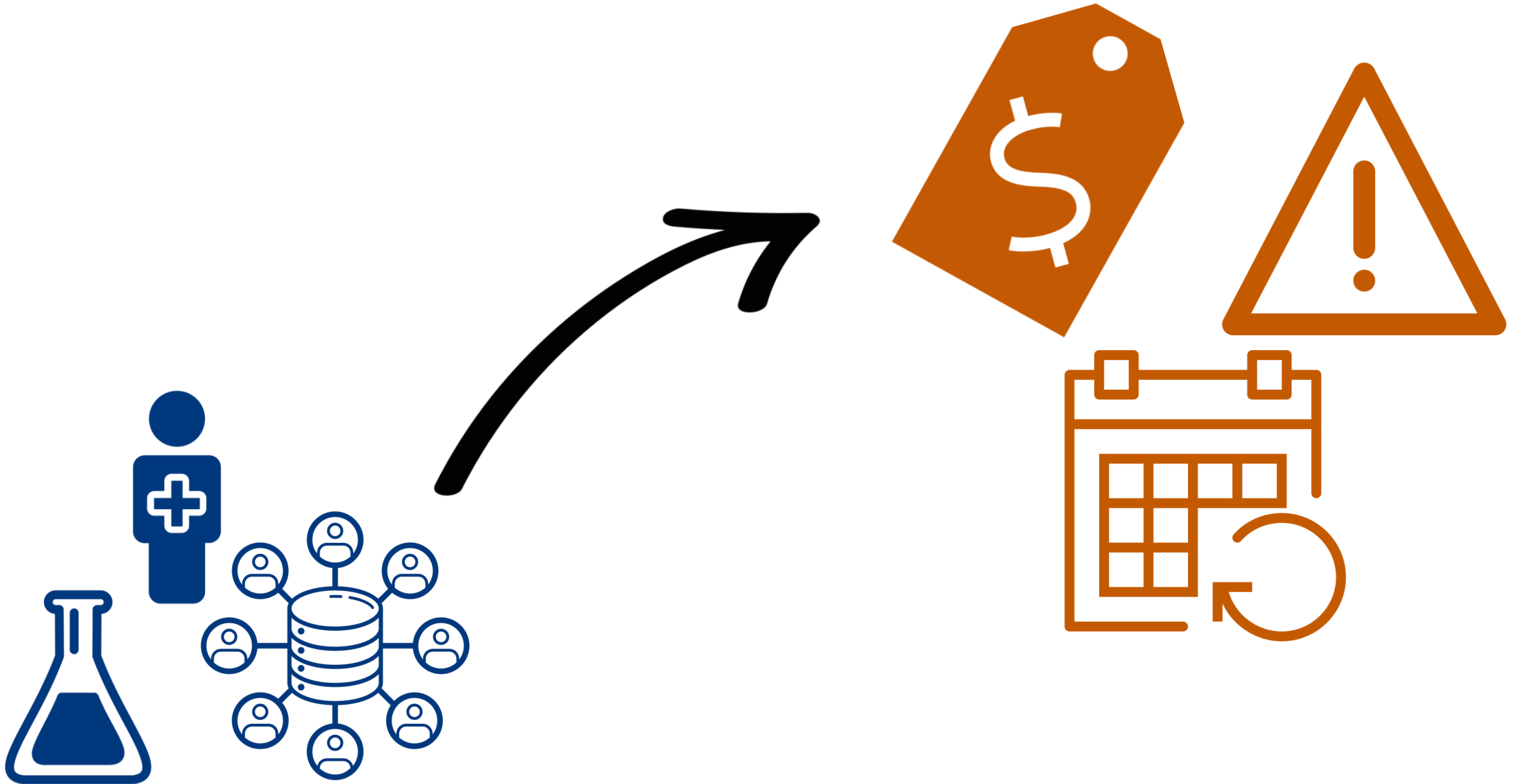
**Study  
planning**

**Study conduct: Recruitment, retention,  
data collection**

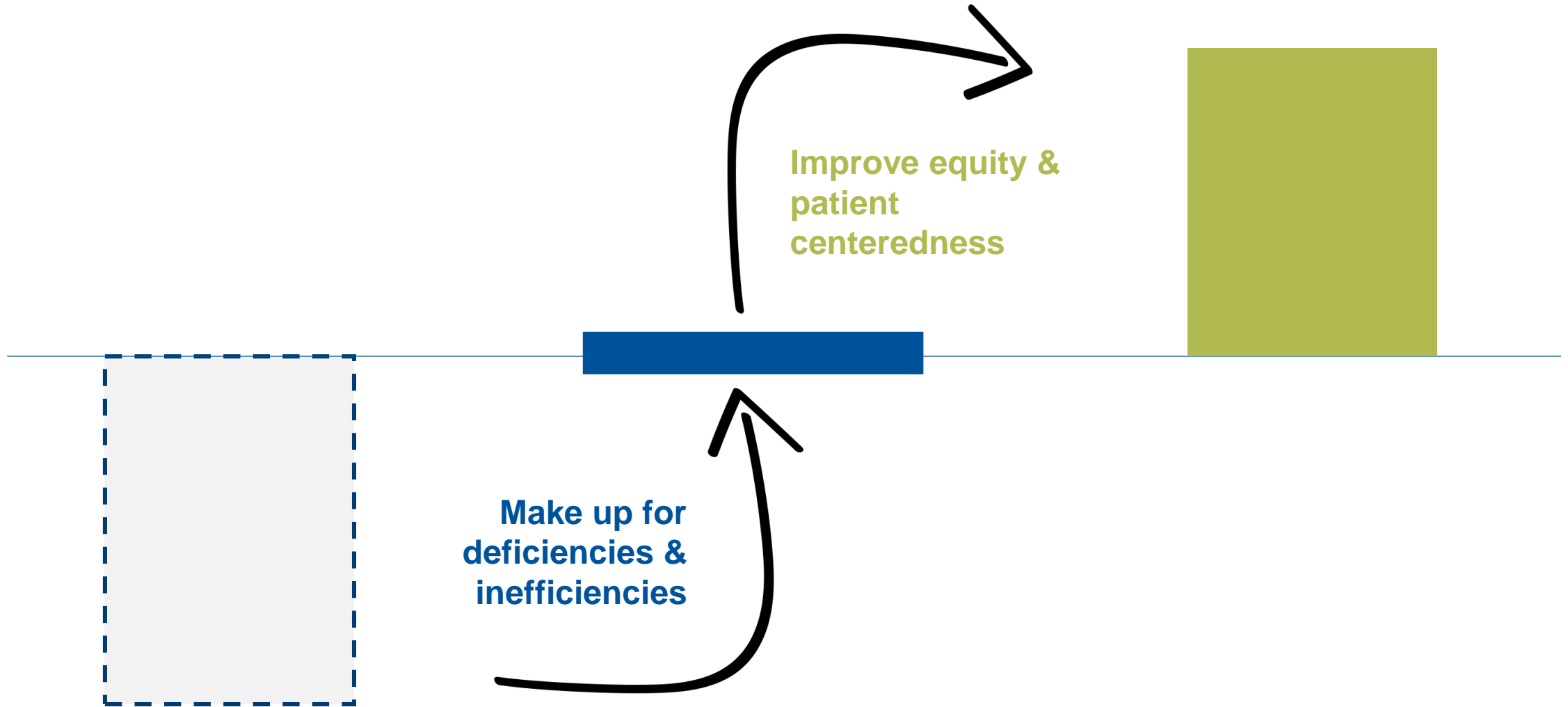
**Data analysis**

**What we'll talk about today**

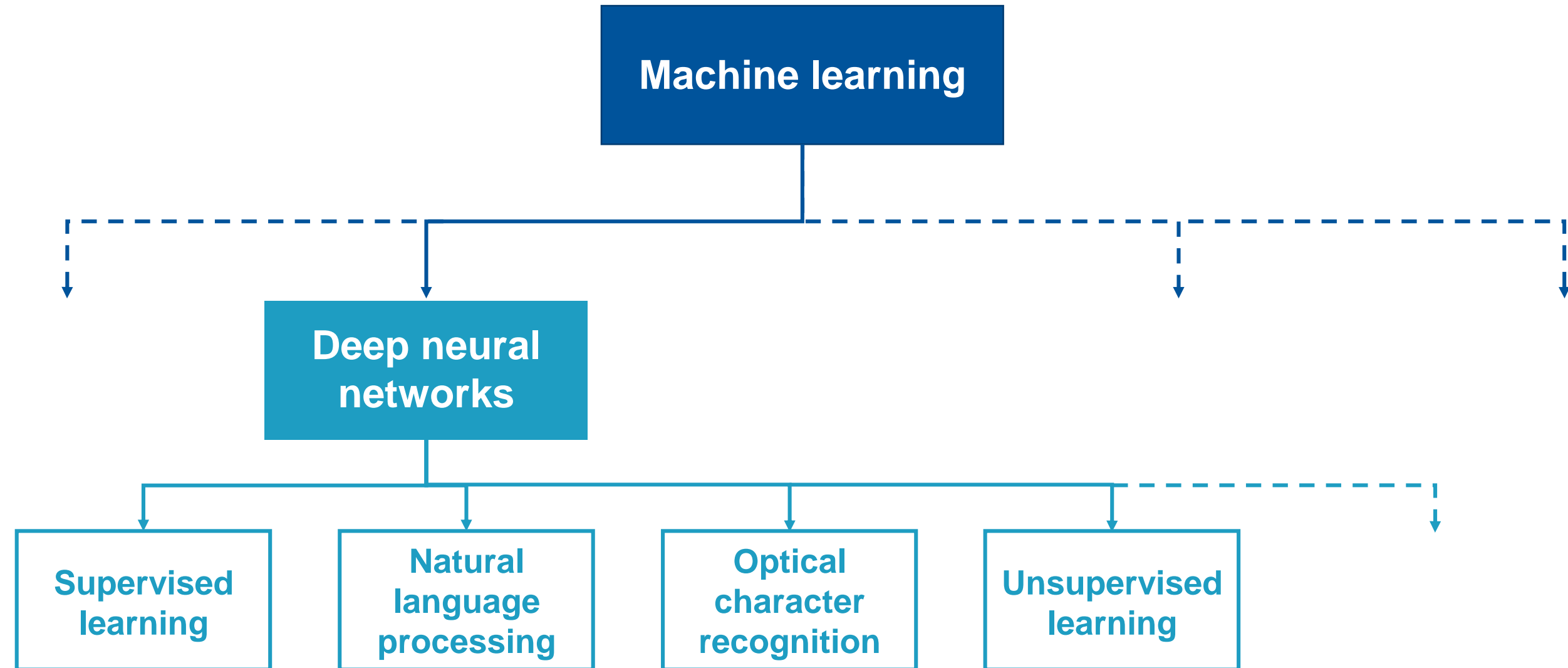
# Past: Opportunities for clinical research improvement



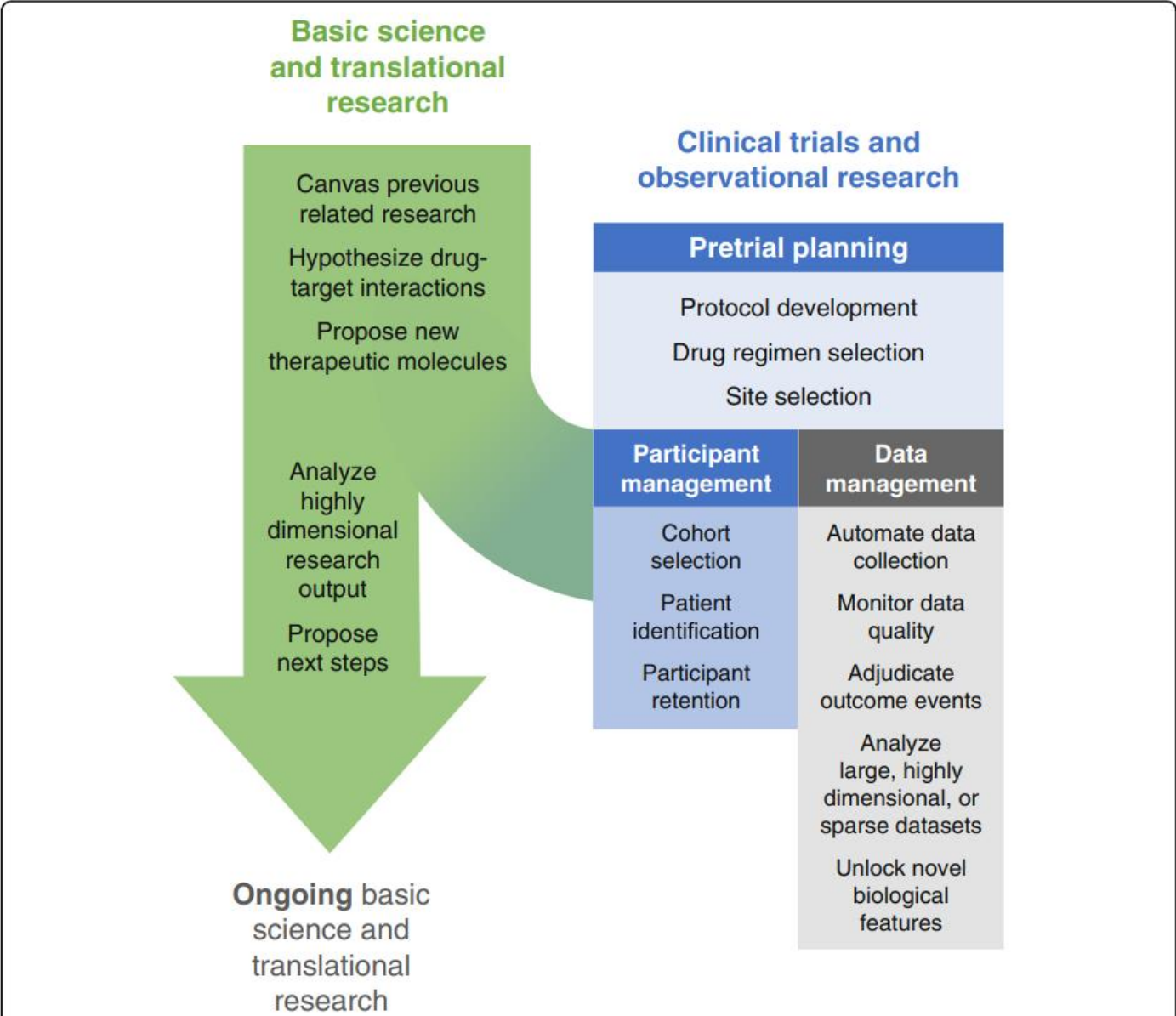
# Past: Opportunities for clinical research improvement



# Present: What can we do with ML in CR?



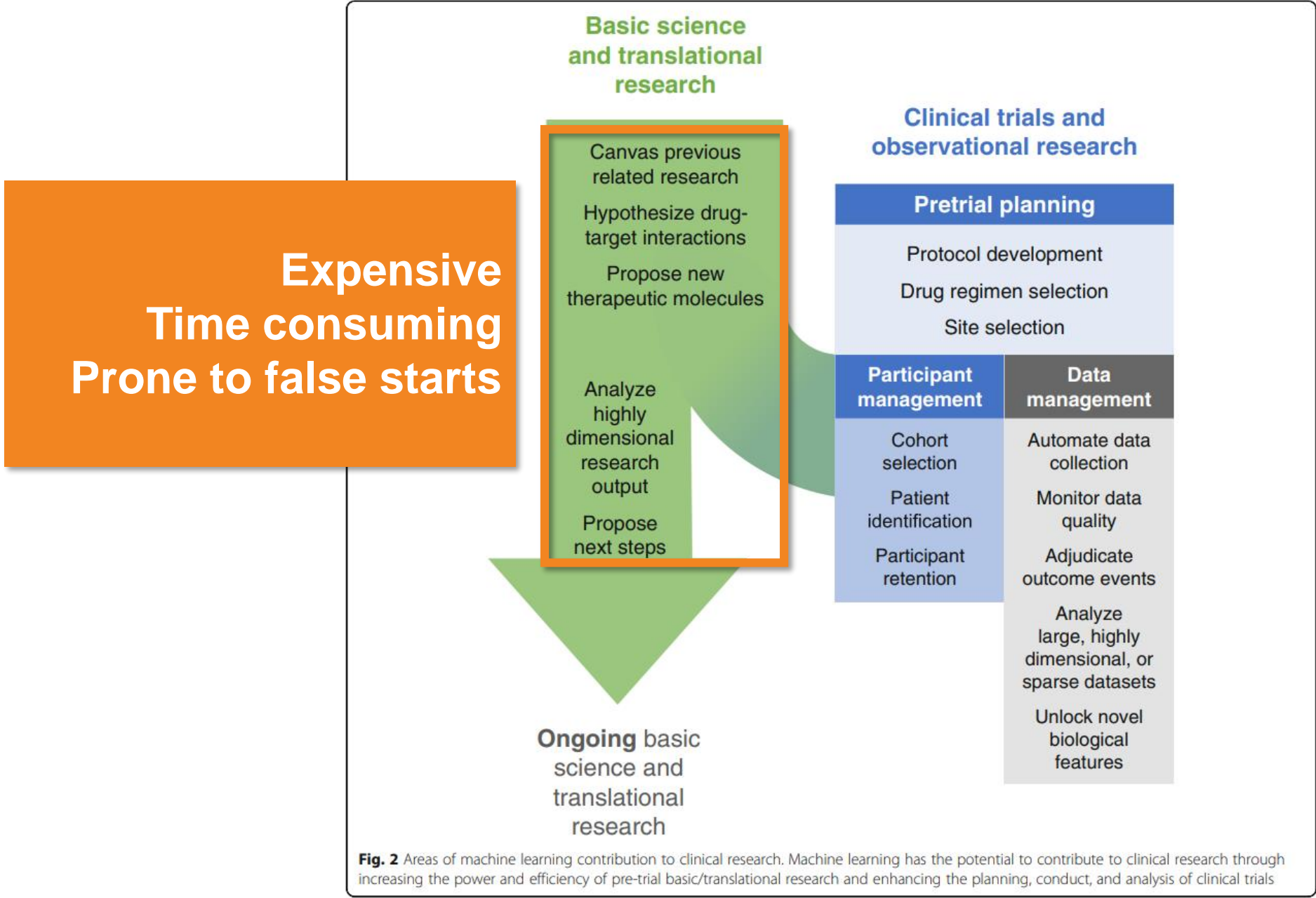
# Present: What can we do with ML in CR?



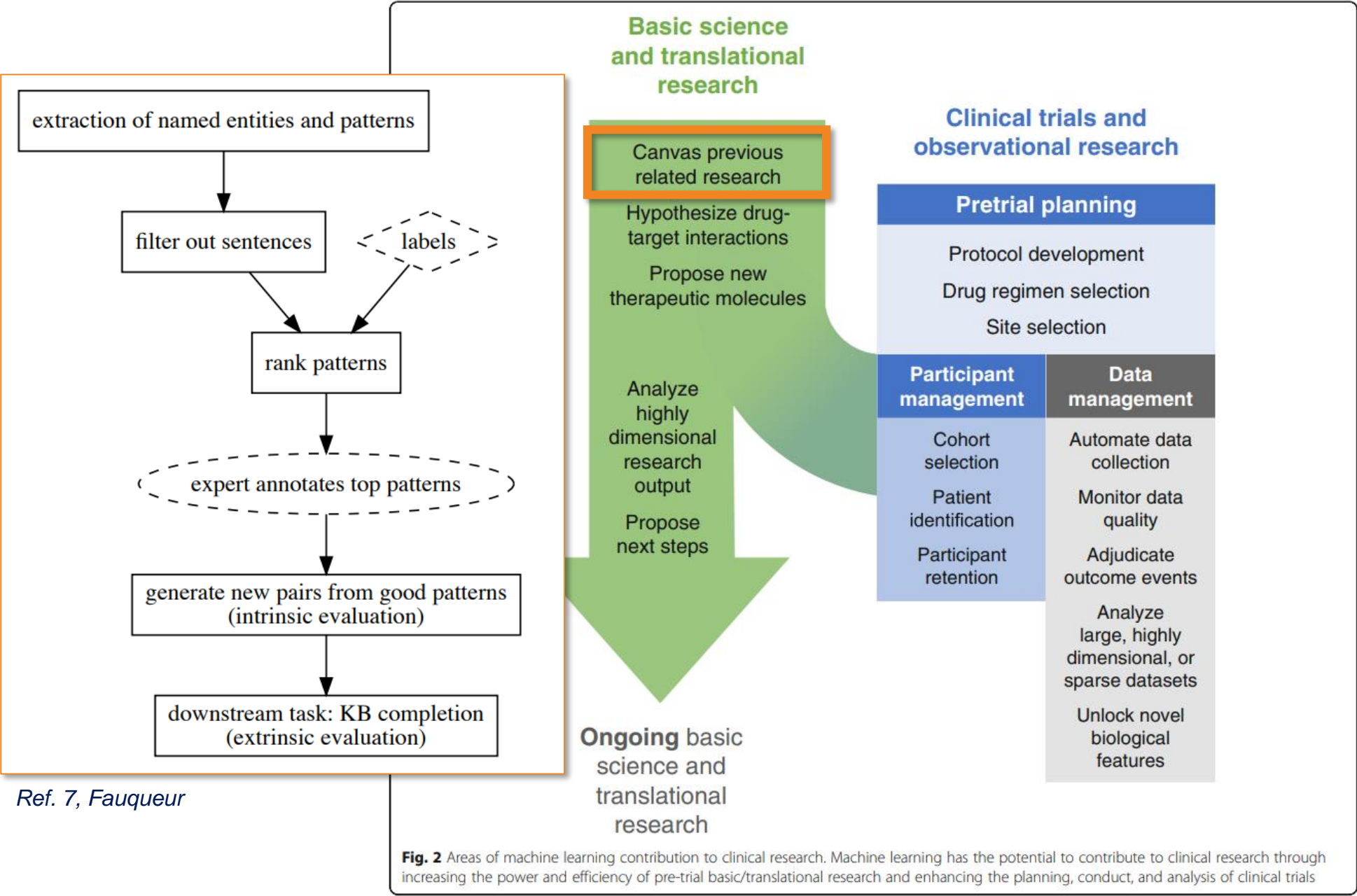
**Fig. 2** Areas of machine learning contribution to clinical research. Machine learning has the potential to contribute to clinical research through increasing the power and efficiency of pre-trial basic/translational research and enhancing the planning, conduct, and analysis of clinical trials



# Present: What can we do with ML in CR?

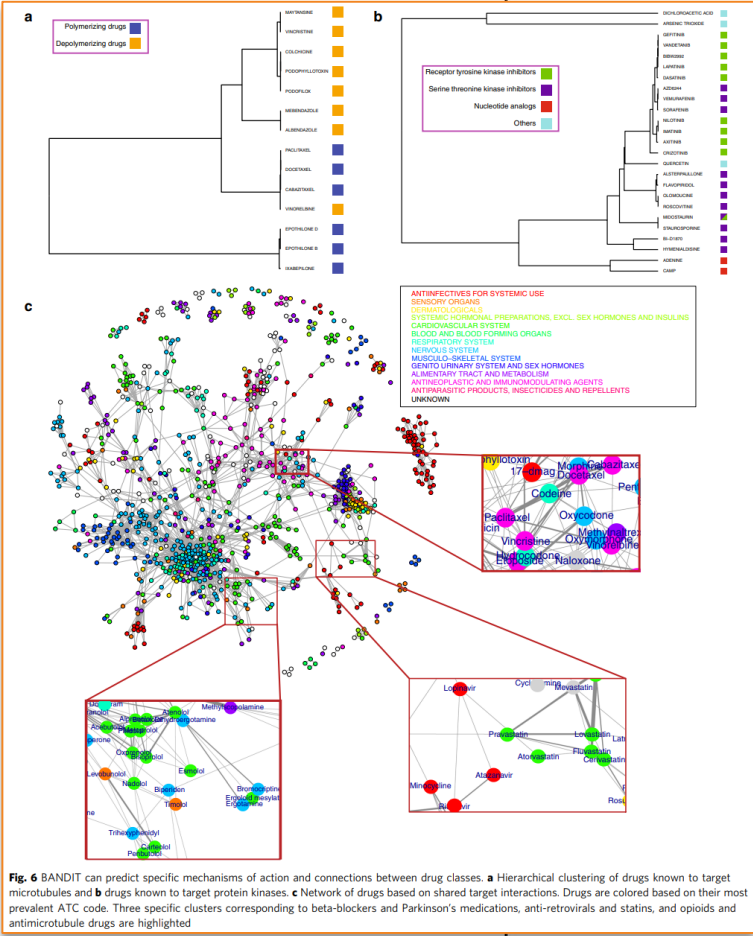


# Present: What can we do with ML in CR?

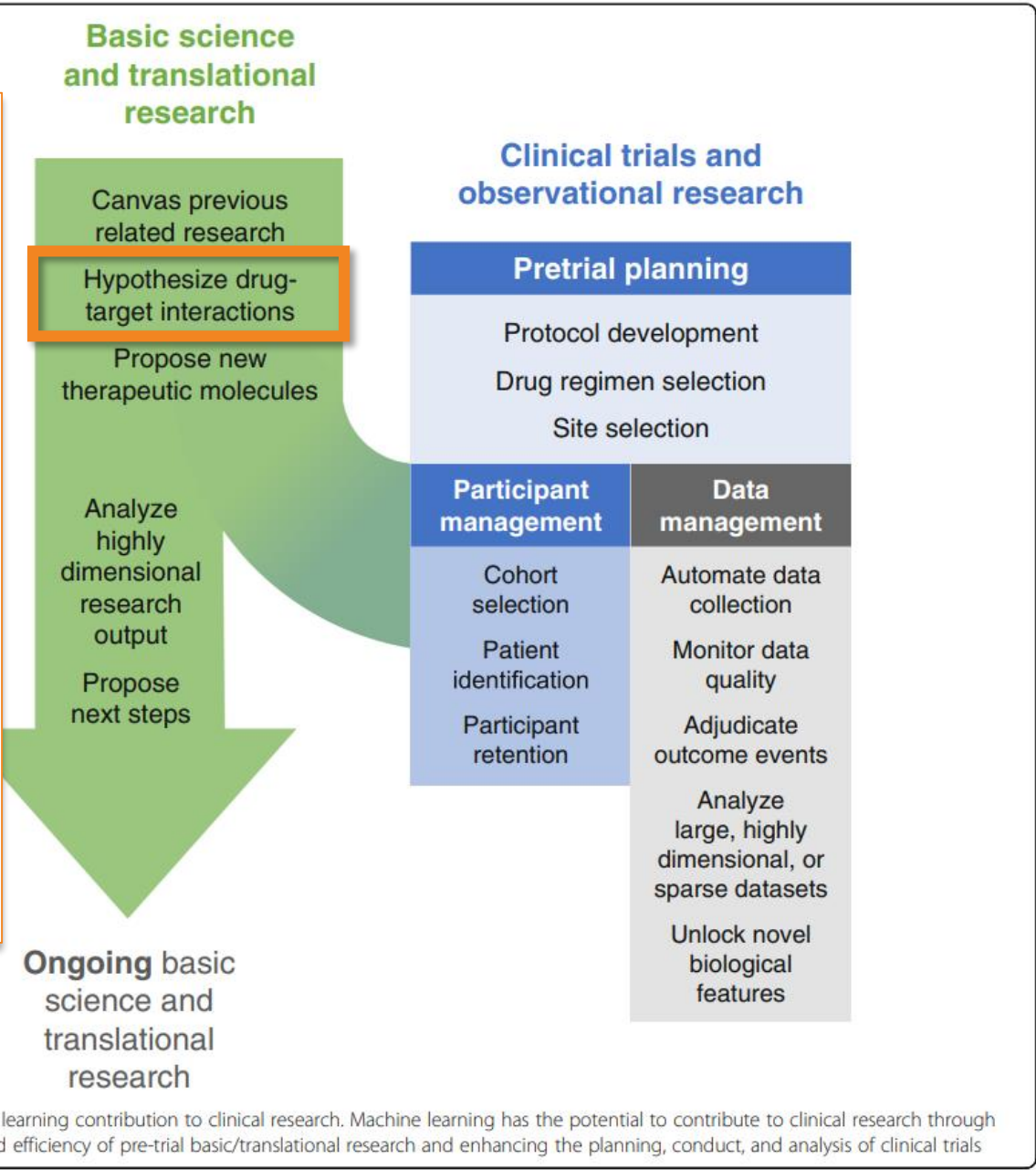


**Fig. 2** Areas of machine learning contribution to clinical research. Machine learning has the potential to contribute to clinical research through increasing the power and efficiency of pre-trial basic/translational research and enhancing the planning, conduct, and analysis of clinical trials

# Present: What can we do with ML in CR?

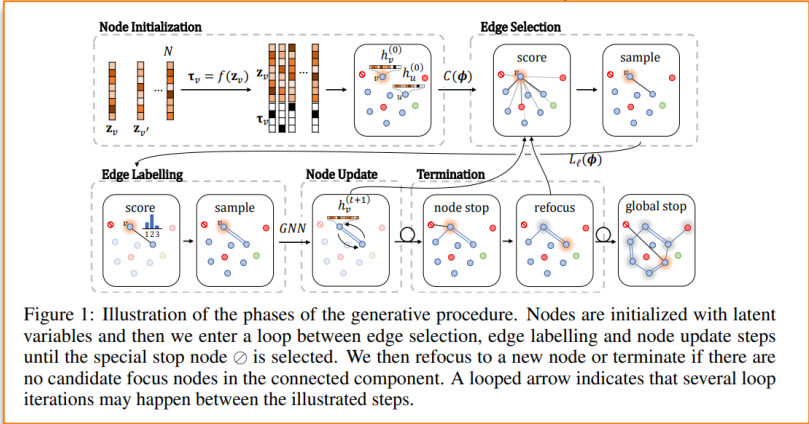


Ref. 8, Madhukar





# Present: What can we do with ML in CR?



Ref. 9, Liu

Basic science  
and translational  
research

Canvas previous  
related research

Hypothesize drug-  
target interactions

Propose new  
therapeutic molecules

Analyze  
highly  
dimensional  
research  
output  
  
Propose  
next steps

Ongoing basic  
science and  
translational  
research

Clinical trials and  
observational research

Pretrial planning

Protocol development

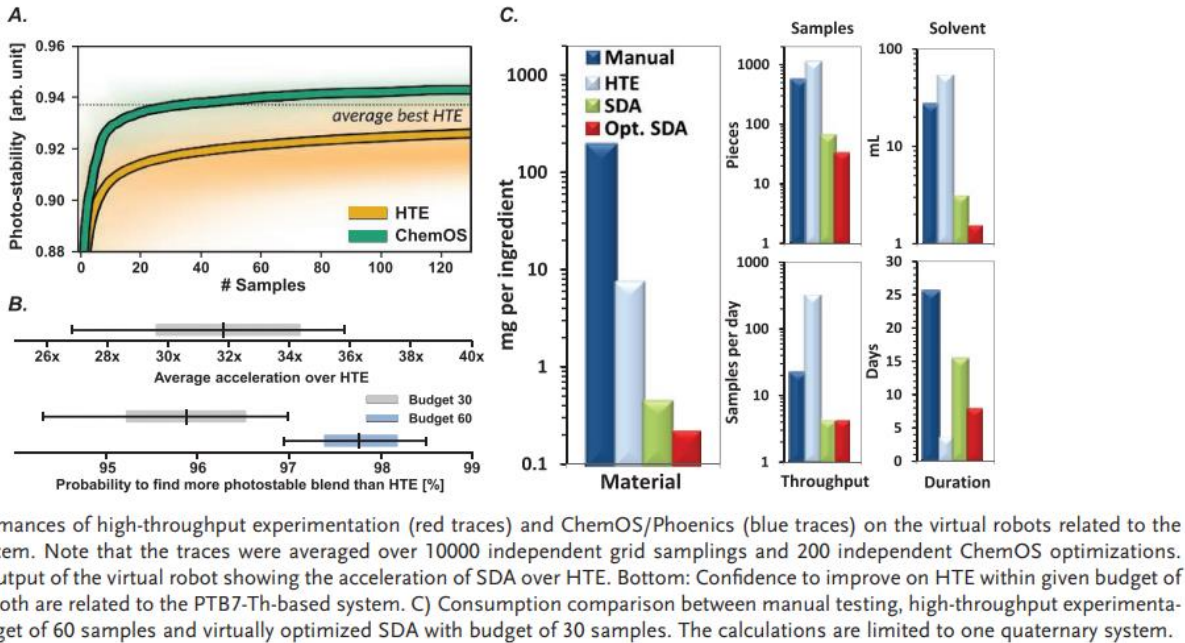
Drug regimen selection

Site selection

Participant

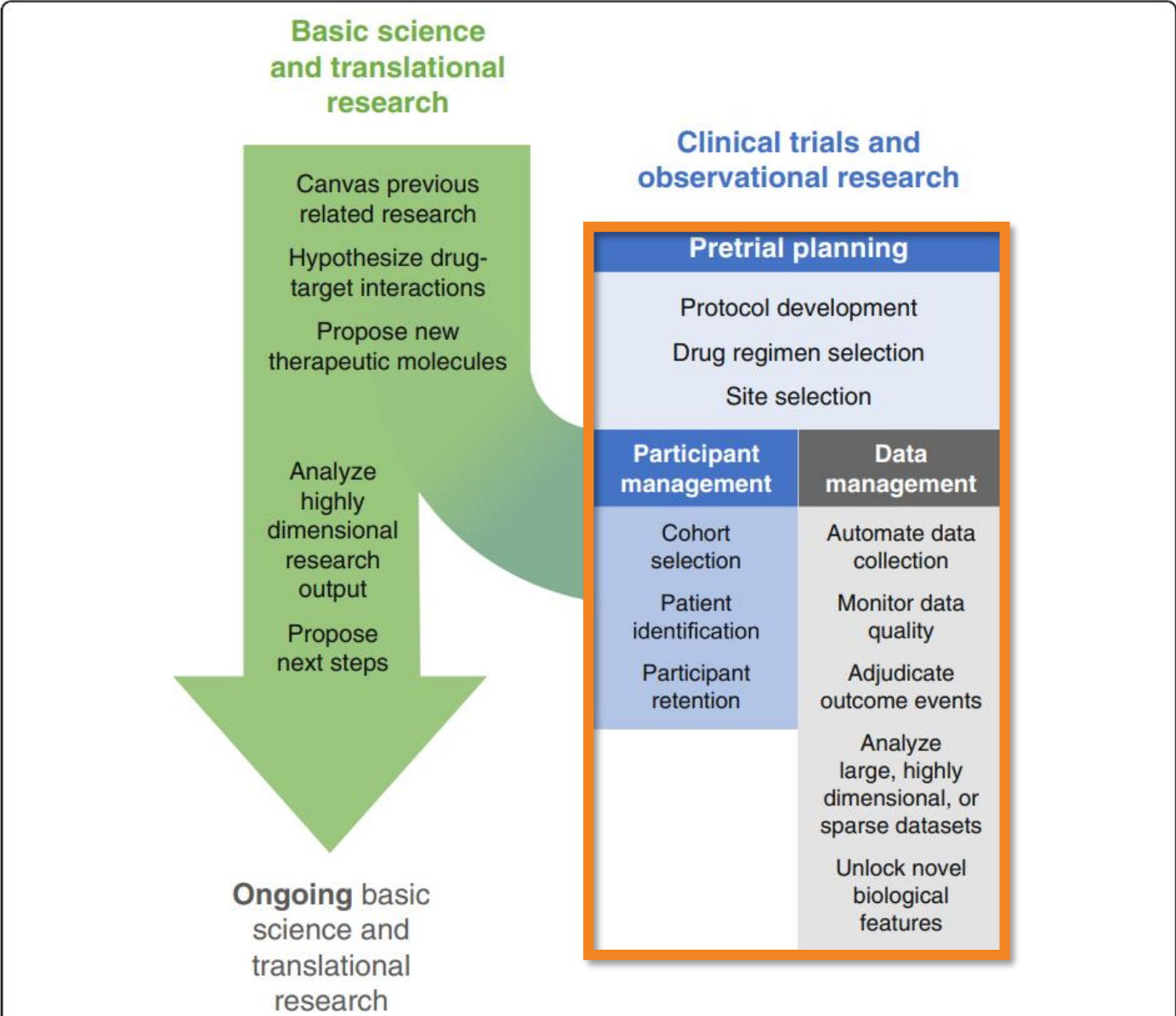
Data

**Fig. 2** Areas of machine learning contribution to clinical research. Machine learning has the potential to contribute to clinical research through increasing the power and efficiency of pre-trial basic/translational research and enhancing the planning, conduct, and analysis of clinical trials



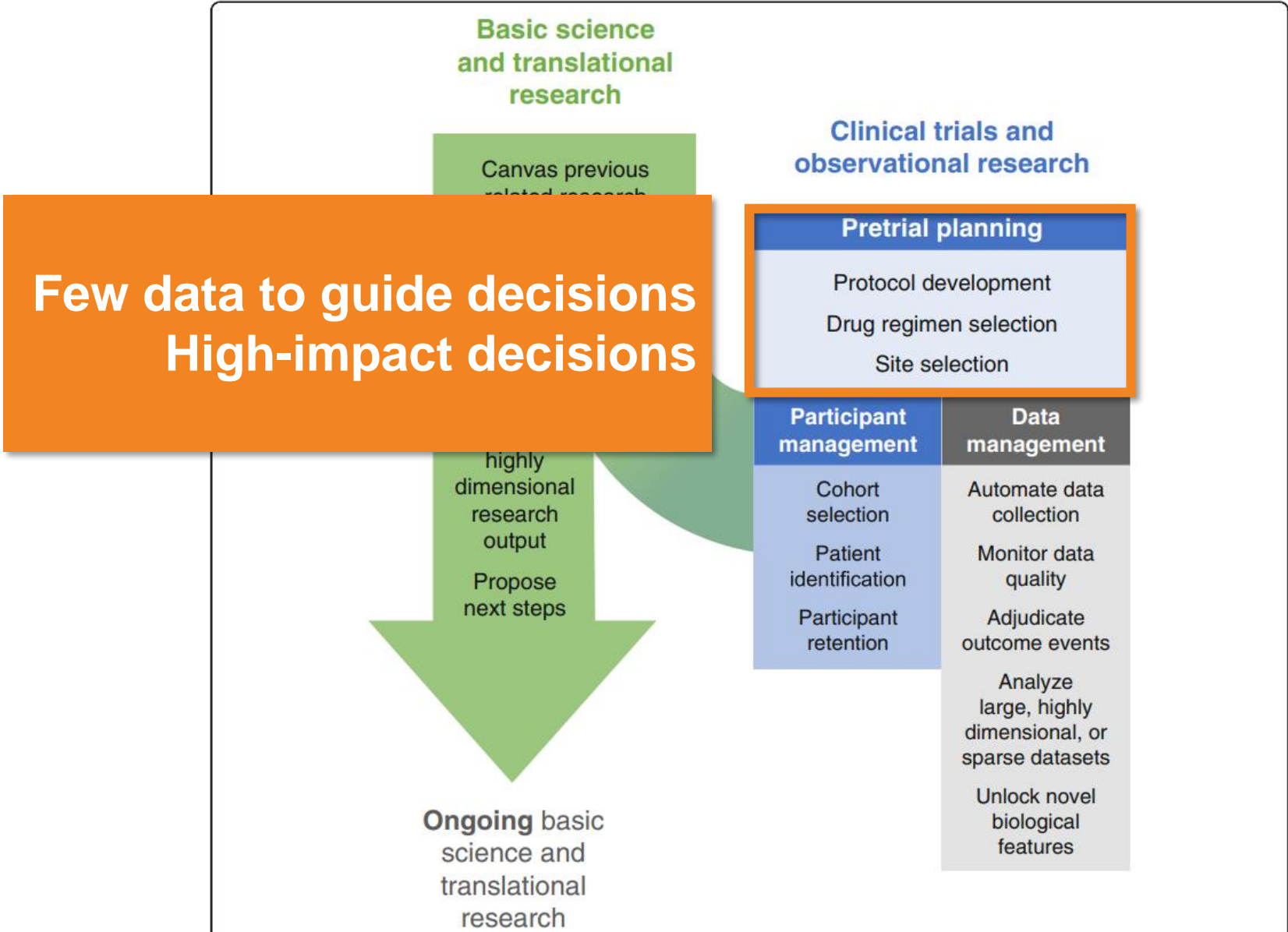
Ref. 10, Langner

# Present: What can we do with ML in CR?



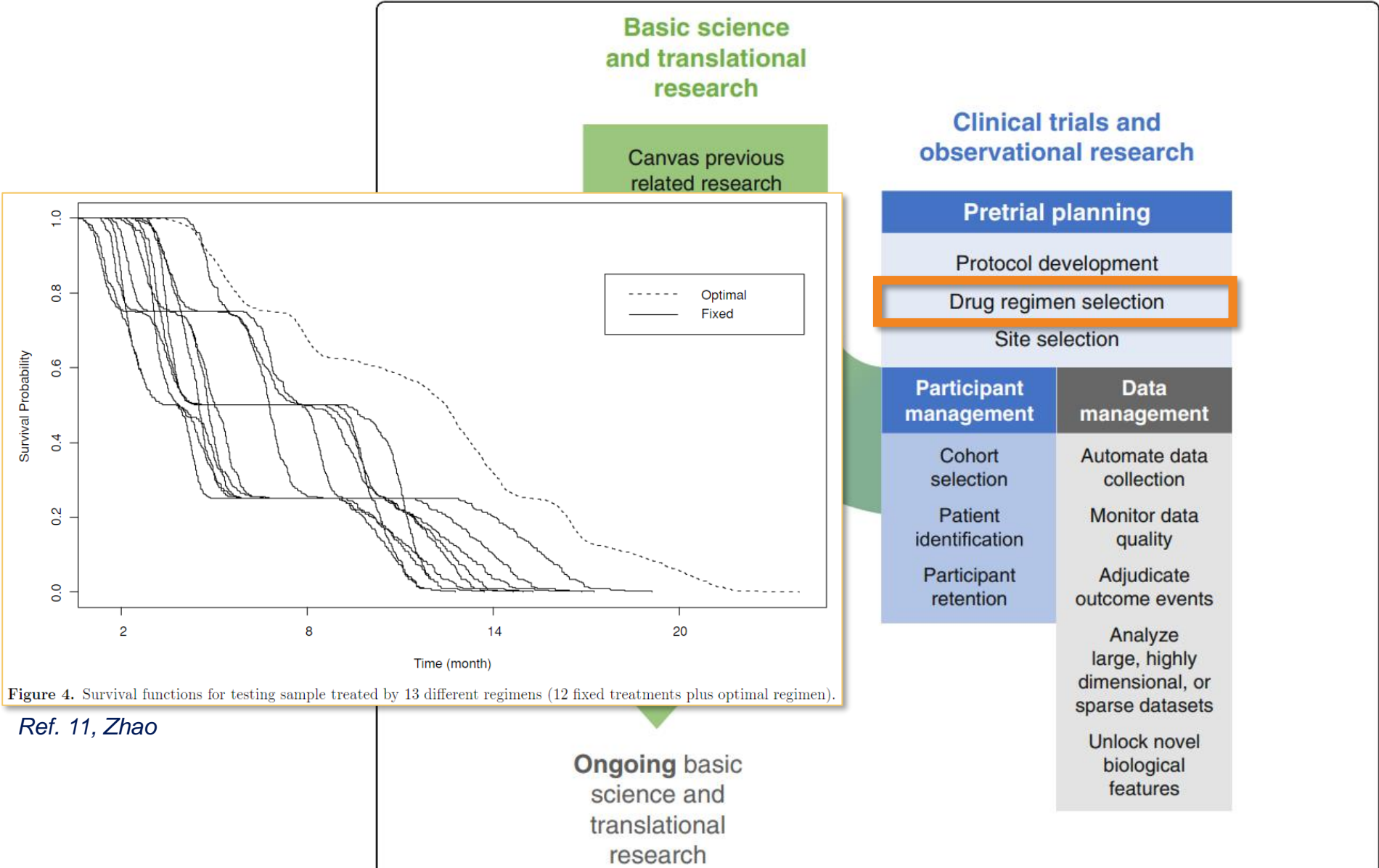
**Fig. 2** Areas of machine learning contribution to clinical research. Machine learning has the potential to contribute to clinical research through increasing the power and efficiency of pre-trial basic/translational research and enhancing the planning, conduct, and analysis of clinical trials

# Present: What can we do with ML in CR?



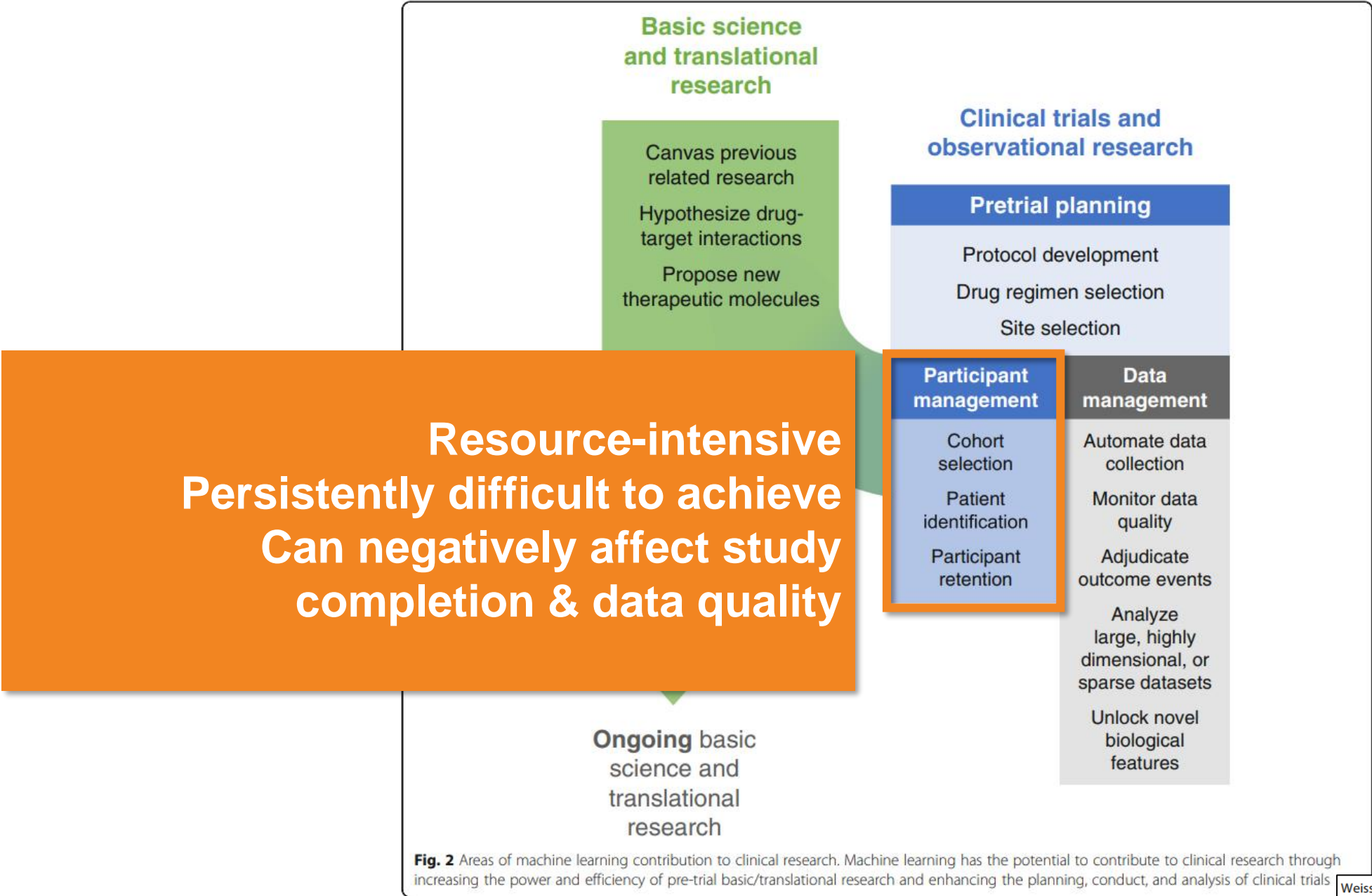
**Fig. 2** Areas of machine learning contribution to clinical research. Machine learning has the potential to contribute to clinical research through increasing the power and efficiency of pre-trial basic/translational research and enhancing the planning, conduct, and analysis of clinical trials

# Present: What can we do with ML in CR?





# Present: What can we do with ML in CR?



**Fig. 2** Areas of machine learning contribution to clinical research. Machine learning has the potential to contribute to clinical research through increasing the power and efficiency of pre-trial basic/translational research and enhancing the planning, conduct, and analysis of clinical trials



# Present: What can we do with ML in CR?

Ref. 12, Seymour

Basic science

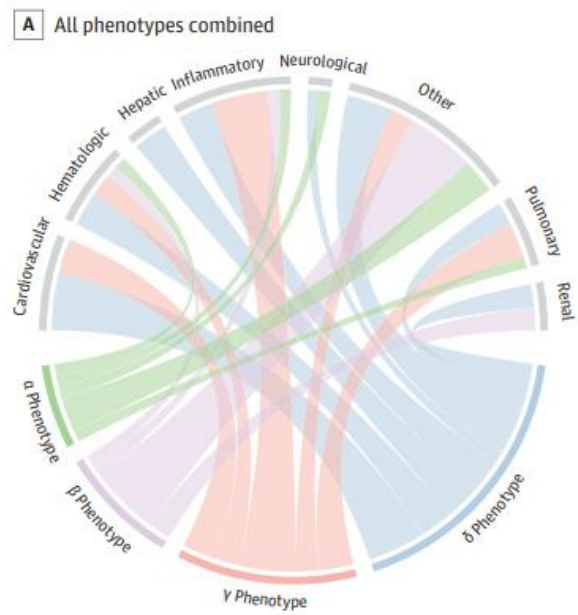
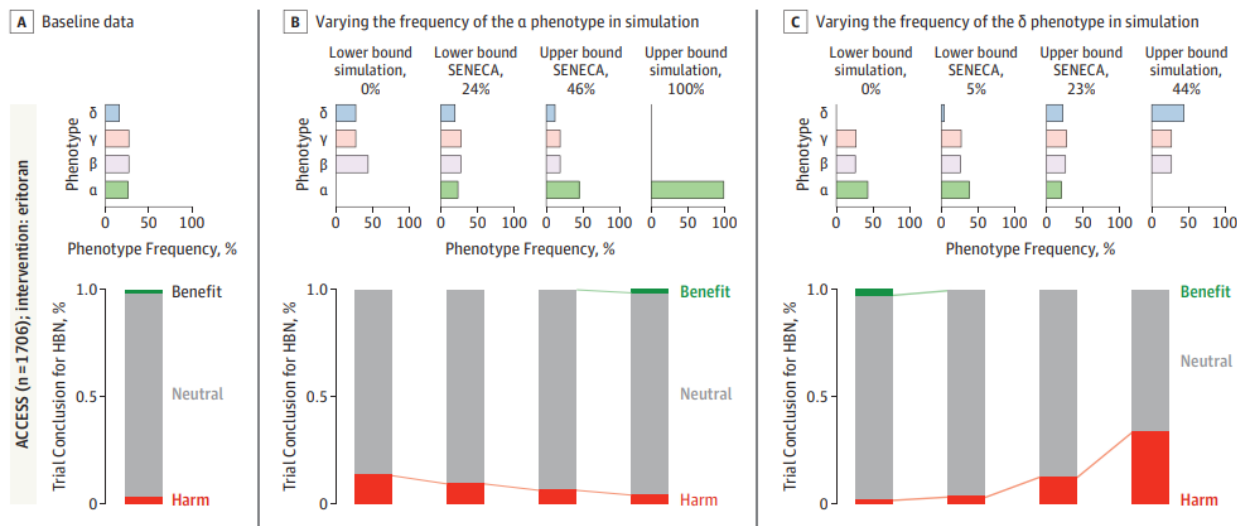


Figure 6. Sensitivity of Clinical Trial Results to the Relative Frequency of Phenotypes in Monte Carlo Simulation



## Clinical trials and observational research

### Pretrial planning

Protocol development  
Drug regimen selection  
Site selection

### Participant management

Cohort selection

Patient identification

Participant retention

### Data management

Automate data collection

Monitor data quality

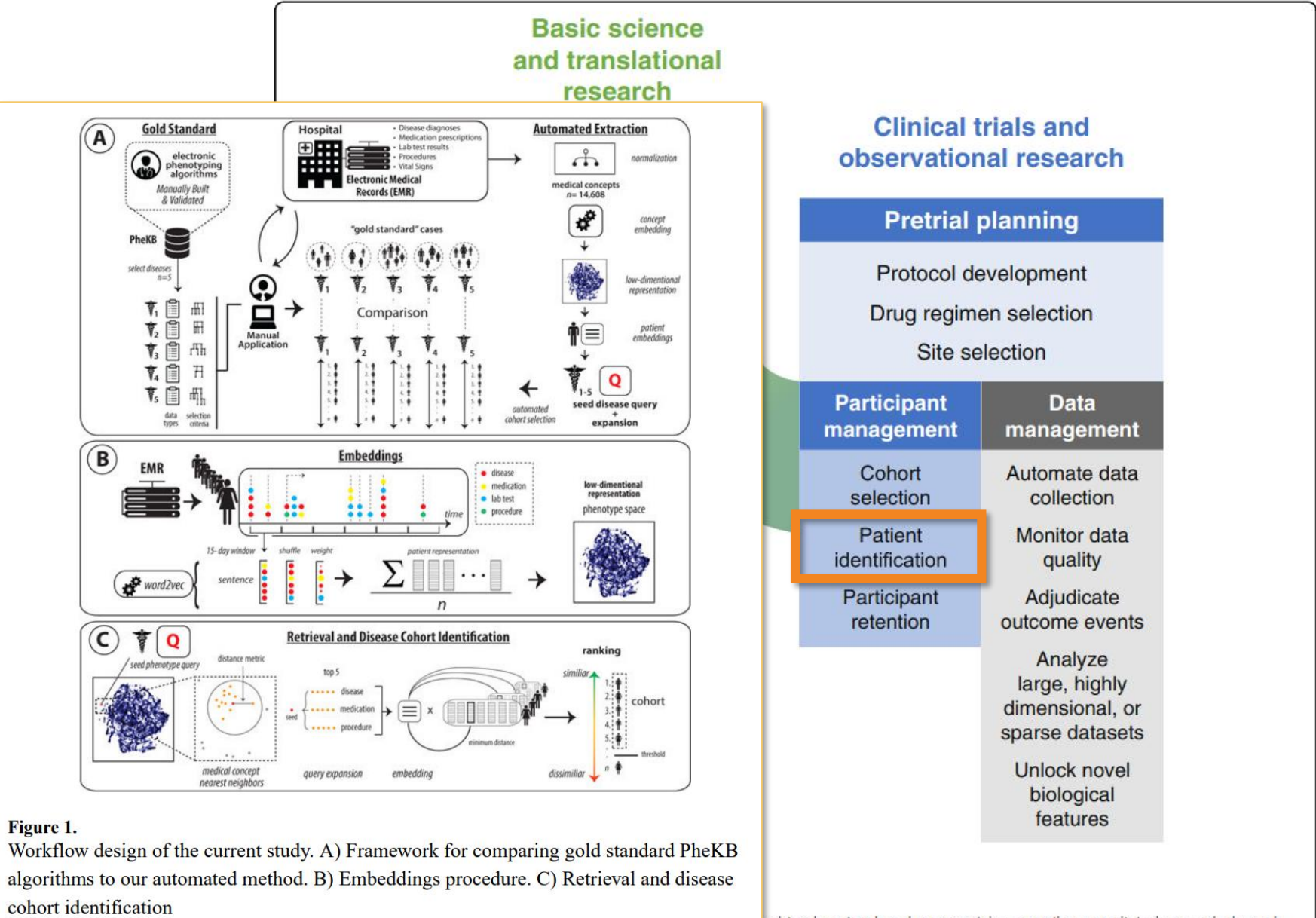
Adjudicate outcome events

Analyze large, highly dimensional, or sparse datasets

Unlock novel biological features

Machine learning has the potential to contribute to clinical research through enhancing the planning, conduct, and analysis of clinical trials

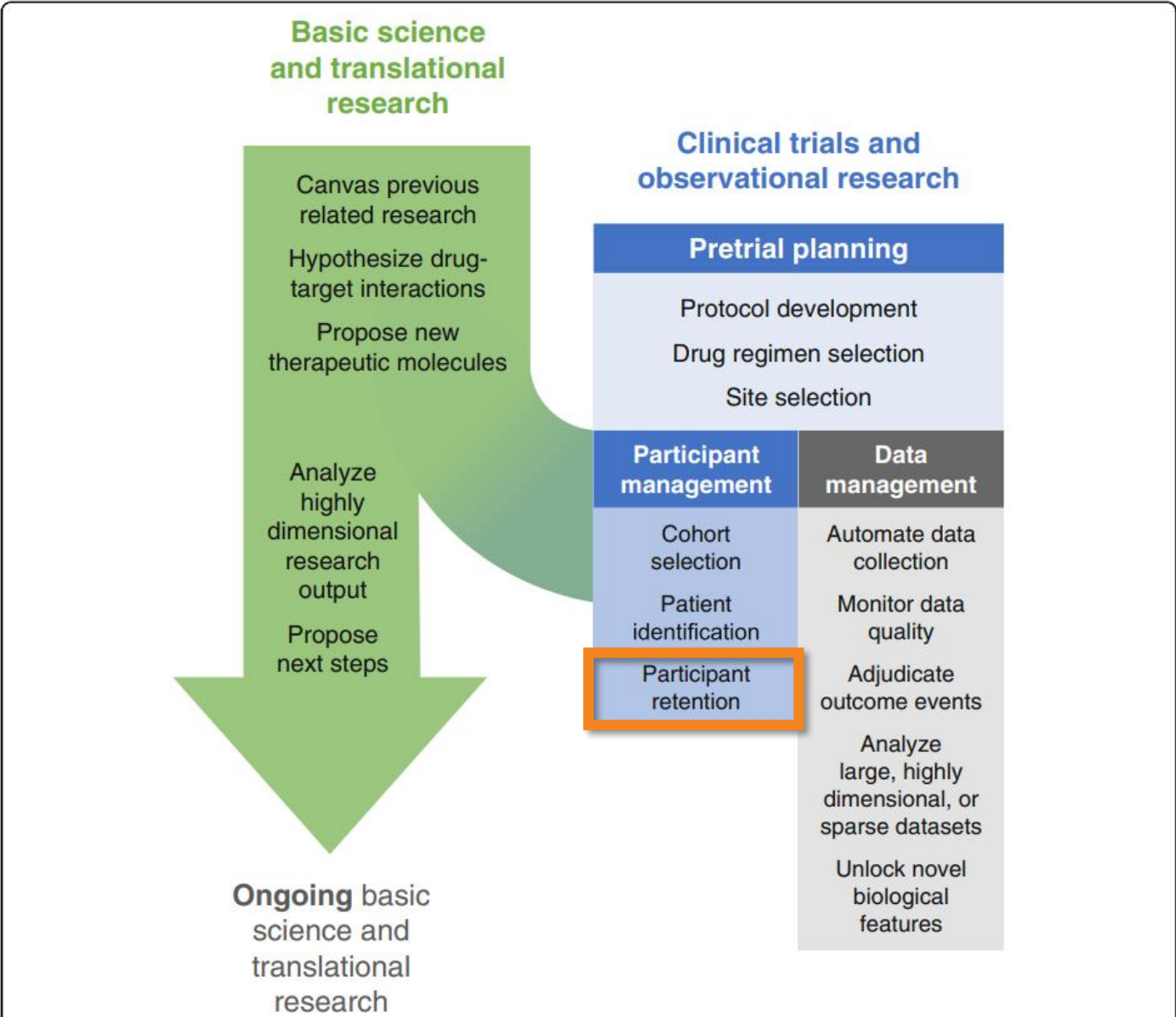
# Present: What can we do with ML in CR?



Ref. 13, Glicksberg

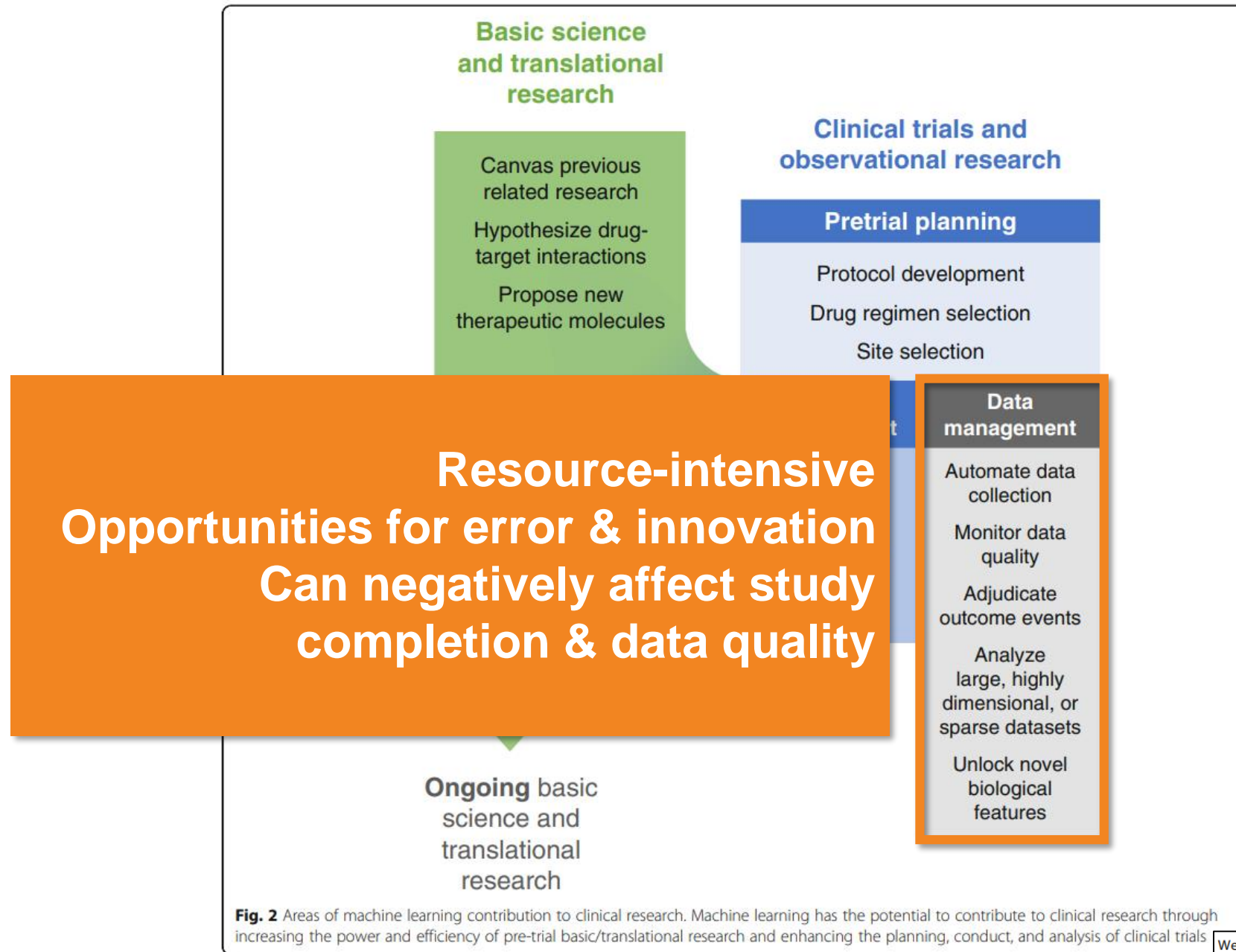
**Fig. 2** Areas of machine learning contribution to clinical research. Machine learning has the potential to contribute to clinical research through increasing the power and efficiency of pre-trial basic/translational research and enhancing the planning, conduct, and analysis of clinical trials

# Present: What can we do with ML in CR?



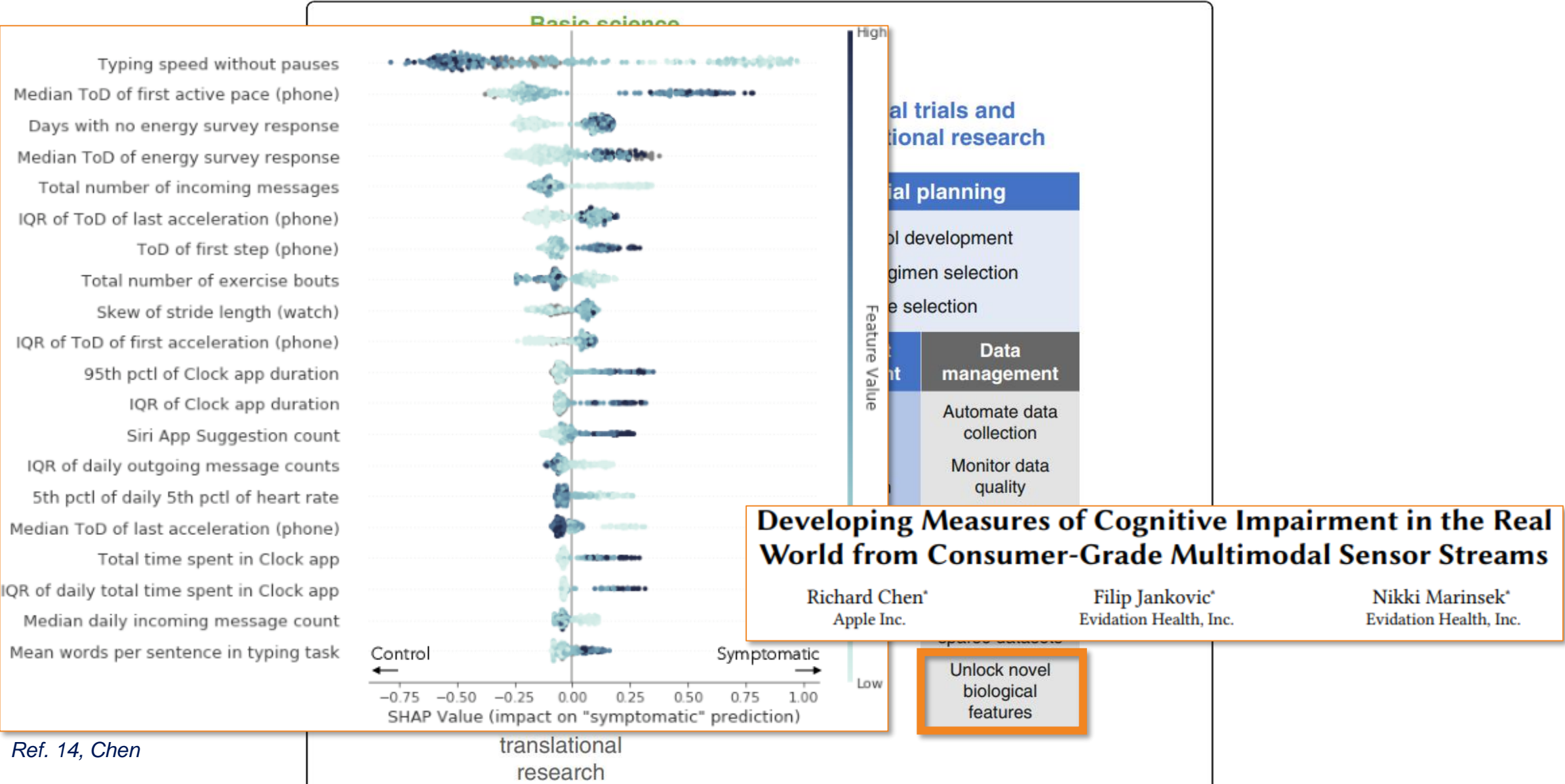
**Fig. 2** Areas of machine learning contribution to clinical research. Machine learning has the potential to contribute to clinical research through increasing the power and efficiency of pre-trial basic/translational research and enhancing the planning, conduct, and analysis of clinical trials

# Present: What can we do with ML in CR?





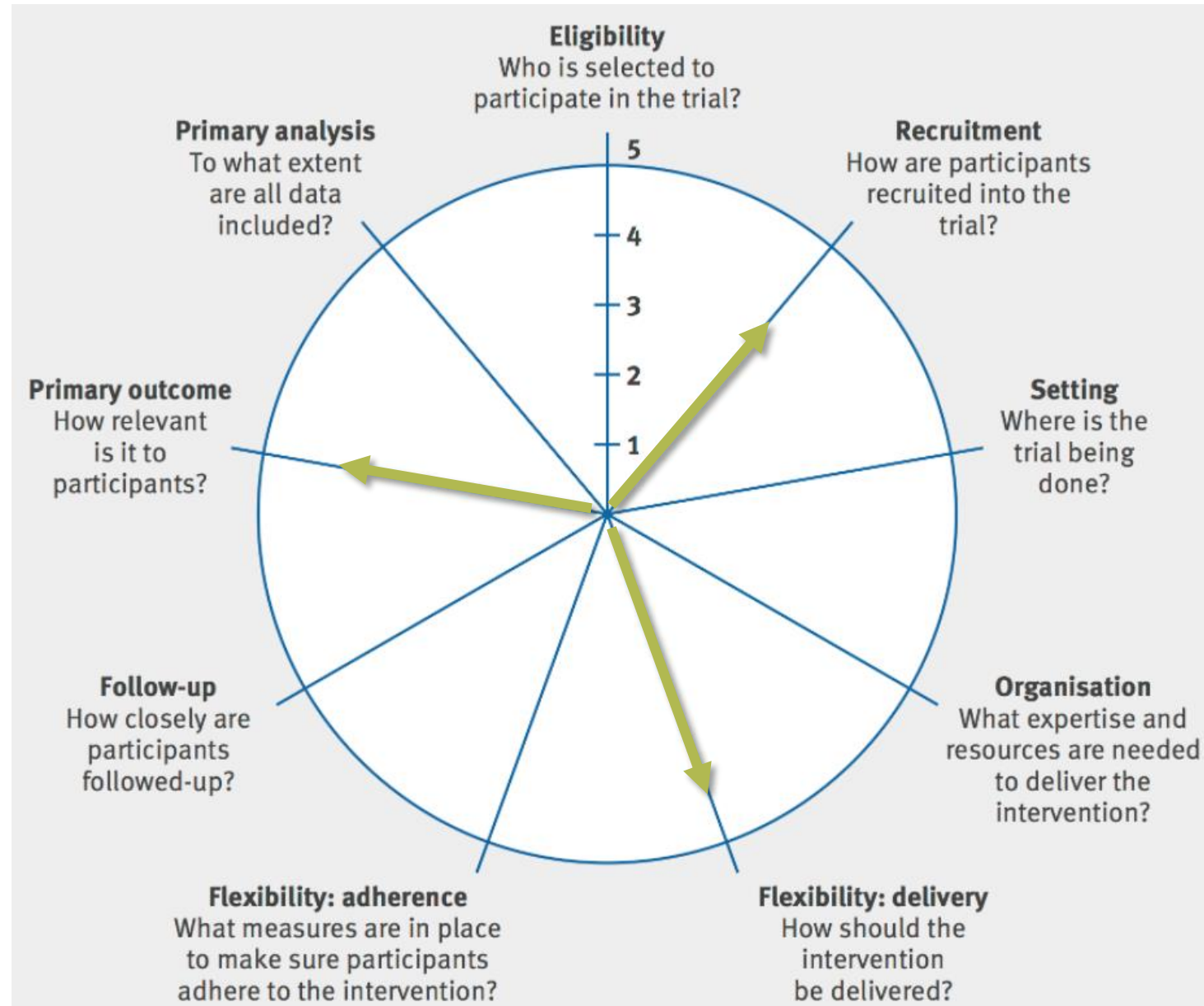
# Present: What can we do with ML in CR?



Ref. 14, Chen

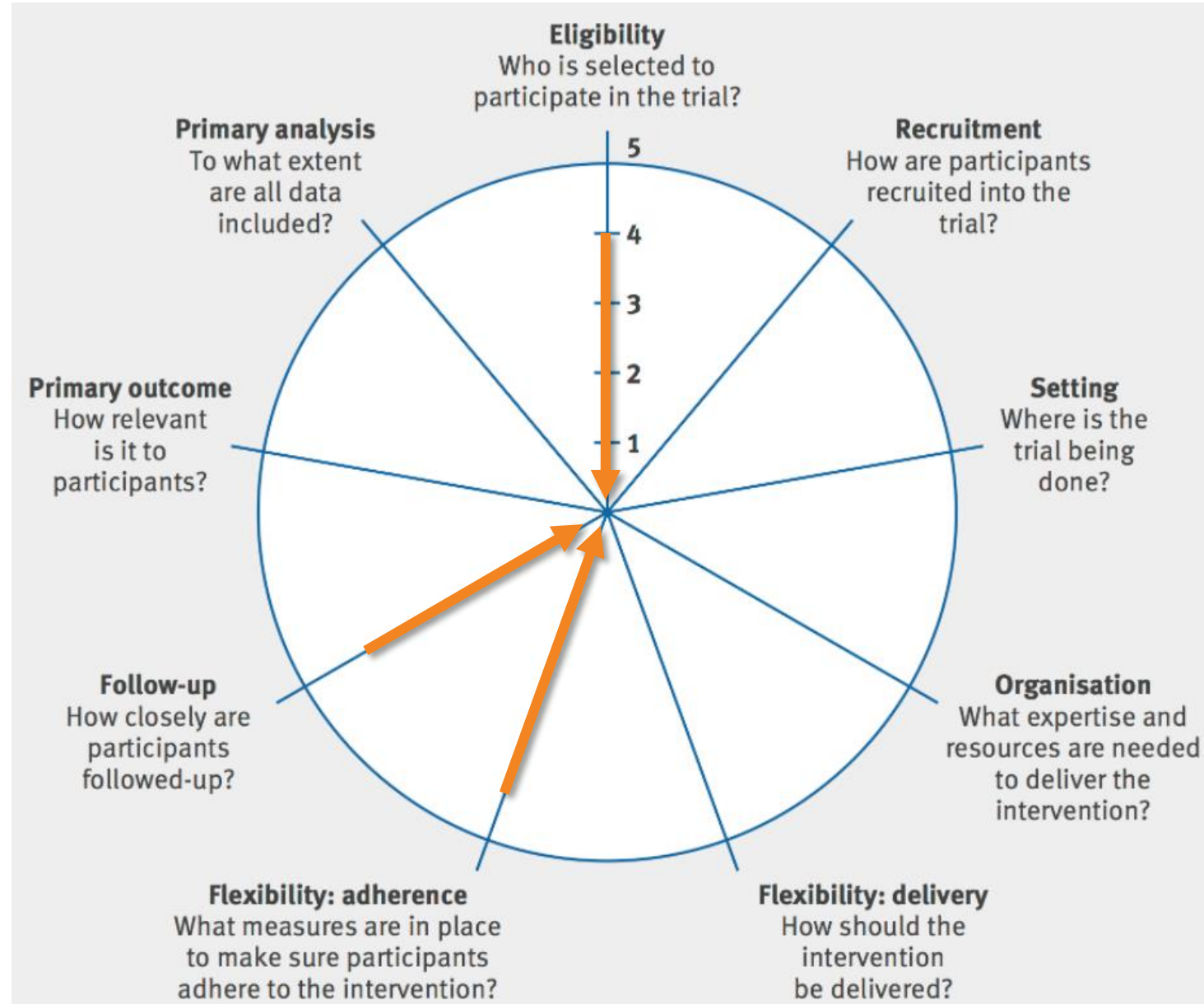
**Fig. 2** Areas of machine learning contribution to clinical research. Machine learning has the potential to contribute to clinical research through increasing the power and efficiency of pre-trial basic/translational research and enhancing the planning, conduct, and analysis of clinical trials

# Present: What can we do with ML in CR?

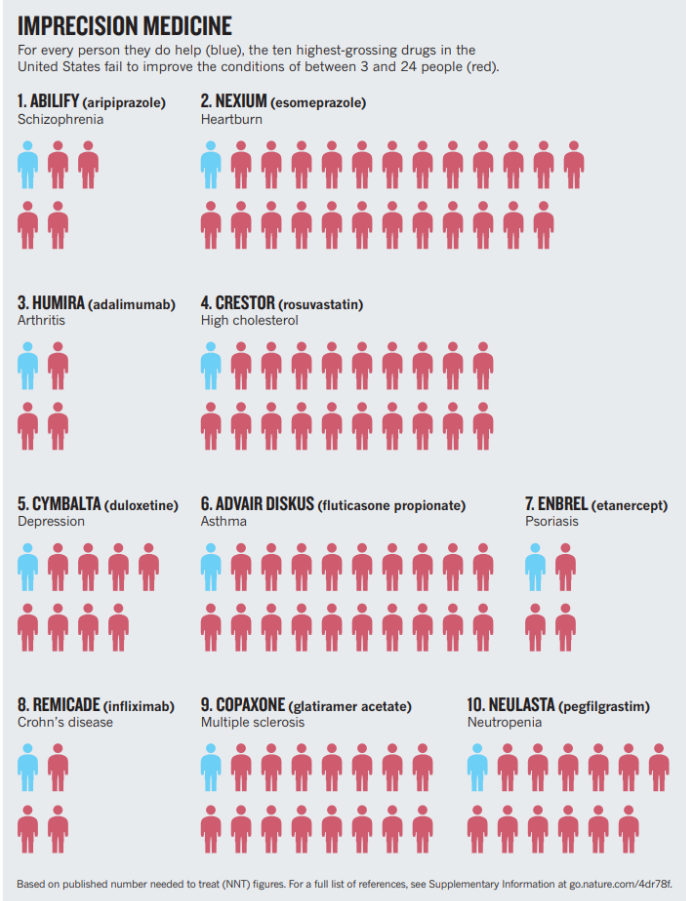


*Ref. 14, Loudon*

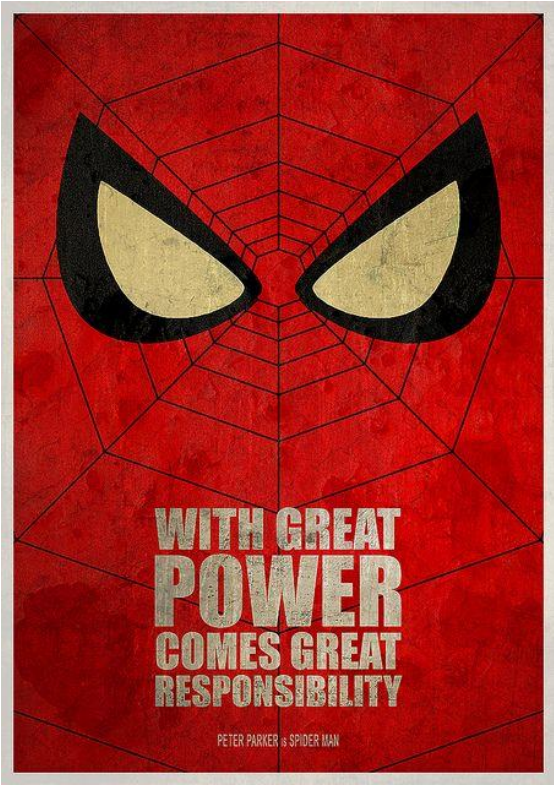
# Present: What can we do with ML in CR?



# Present: What's stopping us from using ML in CR?



Ref. 16, Schork



Ethical risk

Ethical imperative



# Present: What's stopping us from using ML in CR?

Defining the Potential of Machine Learning and Artificial Intelligence Approaches	Erich Huang, <i>Duke Clinical Research Institute</i> Marzyeh Ghassemi, <i>University of Toronto</i>	Transforming Data Surveillance During Trial Conduct	Matthew Roe, <i>Duke Clinical Research Institute</i> Zhaoling Meng, <i>Sanofi</i> Ricardo Henao, <i>Duke University</i>
AI/ML Approaches to Enable Healthcare Delivery vs. Clinical Research	Erich Huang, <i>Duke Clinical Research Institute</i>	Optimization of Trial Operational Conduct	Erich Huang, <i>Duke Clinical Research Institute</i> Mohanish Anand, <i>Pfizer</i> Lucas Glass, <i>IQVIA</i> Bram Zuckerman, <i>FDA</i>
AI/ML Framework for Regulated Clinical Research - Opening Thoughts	Marzyeh Ghassemi, <i>University of Toronto</i>		
Biomarker Discoveries and Drug Target Optimization	James Benoit, <i>Harvard University</i> Daniel Freitag, <i>Bayer</i> Paul Slater, <i>Microsoft</i>	– Monitoring of Patient Compliance and Adherence	Yuan Luo, <i>Northwestern University</i>
		– Advanced Risk-Based Monitoring Techniques	Erich Huang, <i>Duke Clinical Research Institute</i>
Cohort Composition/Phenotyping and Patient Identification	Scott Kollins, <i>Duke University</i> Shameer Khader, <i>AstraZeneca</i> Tristan Naumann, <i>Microsoft</i>	– Endpoint Adjudication	Emmette Hutchison, <i>AstraZeneca</i>
		State of the Art Methods Talk	Marzyeh Ghassemi, <i>University of Toronto</i>

What **can** we do?

What **should** we do and how?

Common Access to Data Sets	Yuan Luo, <i>Northwestern University</i> Marzyeh Ghassemi, <i>University of Toronto</i> Jeff Riesmeyer, <i>Eli Lilly</i> Matthew Diamond, <i>FDA</i> Paul Varghese, <i>Verily</i>	Broader Regulatory Perspective on the Applications of AI/ML for Clinical Trials	Khair ElZarrad, <i>FDA</i>
3rd Party Verification and Certification for Algorithms	Yuan Luo, <i>Northwestern University</i> Stephen Browning, <i>FDA</i> Ricardo Henao, <i>Duke University</i>	Aligning Expectations for Methods	Olivier Elemento, <i>Weill Cornell Medicine</i>
Communities, Challenges, Common Focus Areas	Brian Bot, <i>Sage Bionetworks</i> Masahiro Murakami, <i>Eli Lilly</i> Bray Patrick-Lake, <i>Evidation</i>	Data Needs and Requirements to Optimize and Validate AI/ML Algorithms	Robert Ball, <i>FDA</i> Adarsh Subbaswamy, <i>John Hopkins University</i> Eamon Caddigan, <i>Evidation</i>
		Expectations for Validation of Algorithms	Stephen Browning, <i>FDA</i> Brian Bot, <i>Sage Bionetworks</i> Rajesh Ranganath, <i>New York University</i>
		Preferred Quality Metrics	Boris Brodsky, <i>FDA</i> Michael Hughes, <i>Tufts University</i> Philip Sarocco, <i>Cytokinetics</i> Paul Slater, <i>Microsoft</i>

# Present: What's stopping us from using ML in CR?

## Operational barriers

1. Adequately skilled teams
2. Data:
  - Adequate quantity
  - Multiple sources
  - Adequate quality

# Present: What's stopping us from using ML in CR?

## Operational barriers

1. Adequately skilled teams
2. Data:
  - Adequate quantity
  - Multiple sources
  - Adequate quality

## Philosophical barriers

1. Explainability versus trustworthiness
2. Error and bias

# Present: What's stopping us from using ML in CR?

## Operational barriers

1. Adequately skilled teams
2. Data:
  - Adequate quantity
  - Multiple sources
  - Adequate quality



## Validation Reporting



## Philosophical barriers

1. Explainability versus trustworthiness
2. Error and bias

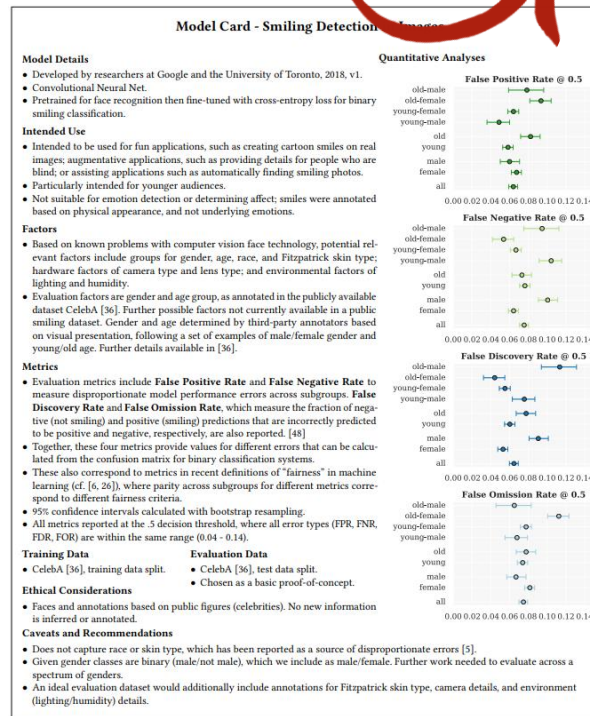


Figure 2: Example Model Card for a smile detector trained and evaluated on the CelebA dataset.

### Model Cards for Model Reporting

Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, Timnit Gebru  
{mmitchellai,simonewu,andrewzaldivar,parkerbarnes,lucyvasserman,benhutch,espitzer,tgebru}@google.com  
deborah.rajai@mail.utoronto.ca

Ref. 17

# Present: What's stopping us from using ML in CR?

trials.ai. Accessed February 2, 2021.

trials.ai

ABOUT USLAB NOTESCONTACT US

Our Smart Protocol system aims to automate and add intelligence to clinical trial design with AI

https://deep6.ai/how-it-works/

DEEP 6 AI


HOW IT WORKSSCHEDULE A DEMOABOUT USCAREERSBLOGCONTACT

FIND MORE PATIENTS IN MINUTES NOT MONTHS

aicure.com

AiCure

EVENTS & NEWSBLOGLEARNOPENDBMPARTNERSHIPSCOMPANYCONTACT



THE RIGHT DOSE FOR THE RIGHT PATIENT

AiCure delivers a compliant, scalable AI platform to maximize the impact of data on clinical research and end-to-end operations, from pre-clinical to commercialization.

AiCure brings AI to the life sciences industry to guide data-driven decision-making for more meaningful clinical trials, optimized drug development, and improved business operations.

Fig. 2  
increas



BULLFROG AI

HOME TECHNOLOGY PARTNERING SOLUTIONS PIPELINE LEADERSHIP INVESTORS MEDIA

DISCOVERING THERAPEUTICS THAT WOULDN'T BE FOUND BEFORE.

88 PERCENT OF DRUGS IN PIPELINE WILL FAIL

8½ YEARS AVERAGE UNTIL READY FOR MARKET

2.6 BILLION AVERAGE COST FOR NEW DRUG

100 BILLION SPENT ON 208 RESEARCH BY TOP 15 PHARMA

204 BILLION TOTAL ANNUAL R&D SPENDING BY 2024; 3% CAGR



Recruit

Increase Enrollment by 50%.  
Prescreen in Under 10 Minutes.

Recruit is an AI-powered solution for clinical sites looking to accelerate patient prescreening,

Mendel.ai

https://www.bullfrogai.com/our-solution/


# Present: What's stopping us from using ML in CR?

Med Health Care and Philos (2016) 19:177–190  
DOI 10.1007/s11019-015-9661-6



## SCIENTIFIC CONTRIBUTION

**“You hoped we would sleep walk into accepting the collection of our data”: controversies surrounding the UK care.data scheme and their wider relevance for biomedical research**



Sigrid Sterckx<sup>1</sup>  • Vojin Rakic<sup>2</sup> • Julian Cockbain<sup>3</sup> • Pascal Borry<sup>4</sup>

*Ref. 18*

## SPECIAL REPORT

**An invisible hand: Patients aren't being told about the AI systems advising their care**



By [Rebecca Robbins](#)  and [Erin Brodwin](#)  July 15, 2020

[Reprints](#)

*Ref. 19*



# **Future: Overcoming barriers to implementation**

**→ Erich & a discussion of data**

# Data Liquidity

*what is it really?*



# Why Data Liquidity?

# Why?

**LEARNING HEALTH** requires unceasing data collection & knowledge generation that is plowed back into patient care

*INTERPRET*

*ANALYZE*

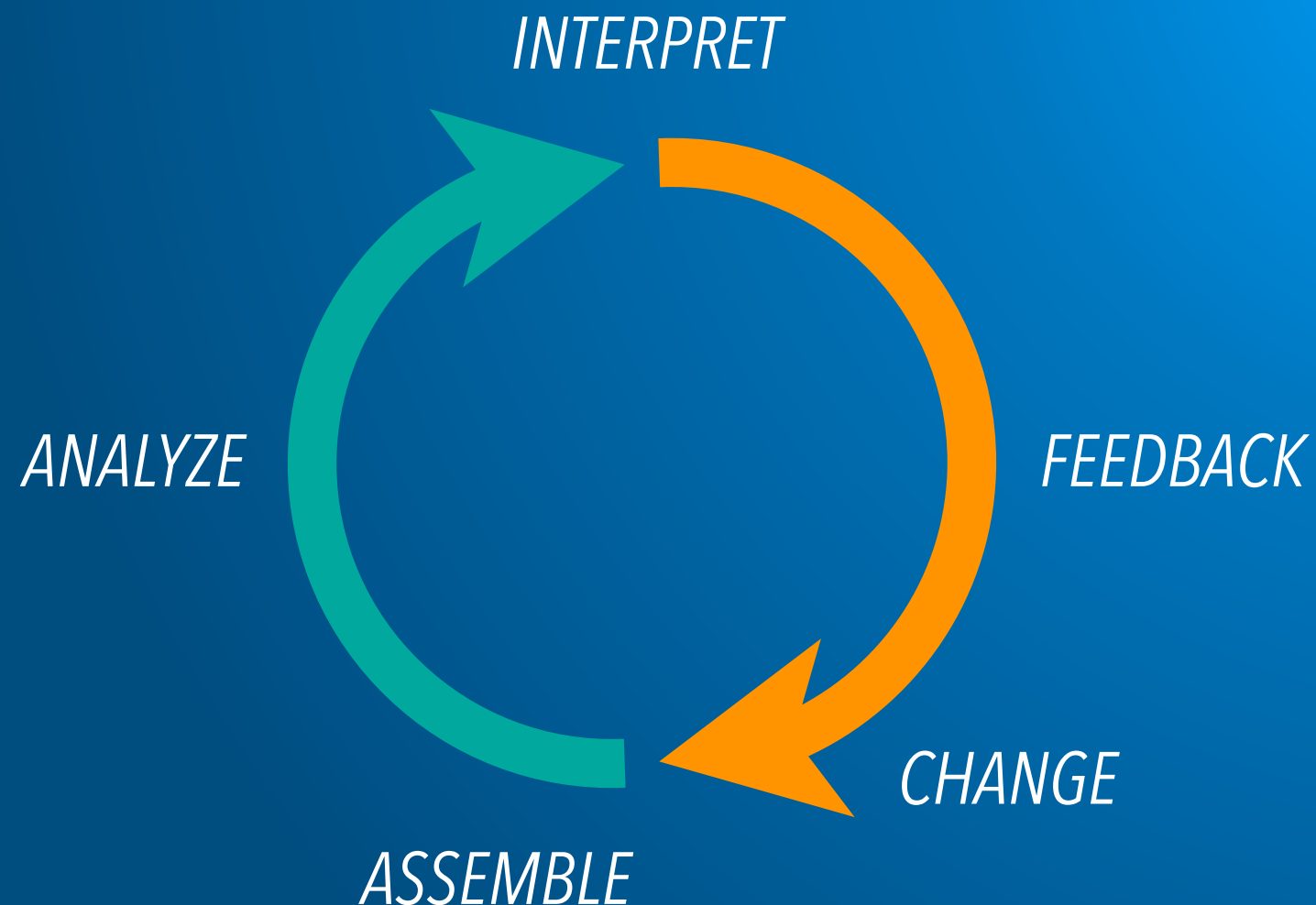
*FEEDBACK*

*CHANGE*

*ASSEMBLE*

# Why?

**LEARNING HEALTH** requires unceasing data collection & knowledge generation that is plowed back into patient care



# A Useful Analogy

# Analogy



# Analogy



*If you can go to an ATM in Antwerp, or anywhere, you can securely access your \$ in your US bank with virtually no friction other than a fee*



# Analogy

Request



Delivery



# Analogy

Request

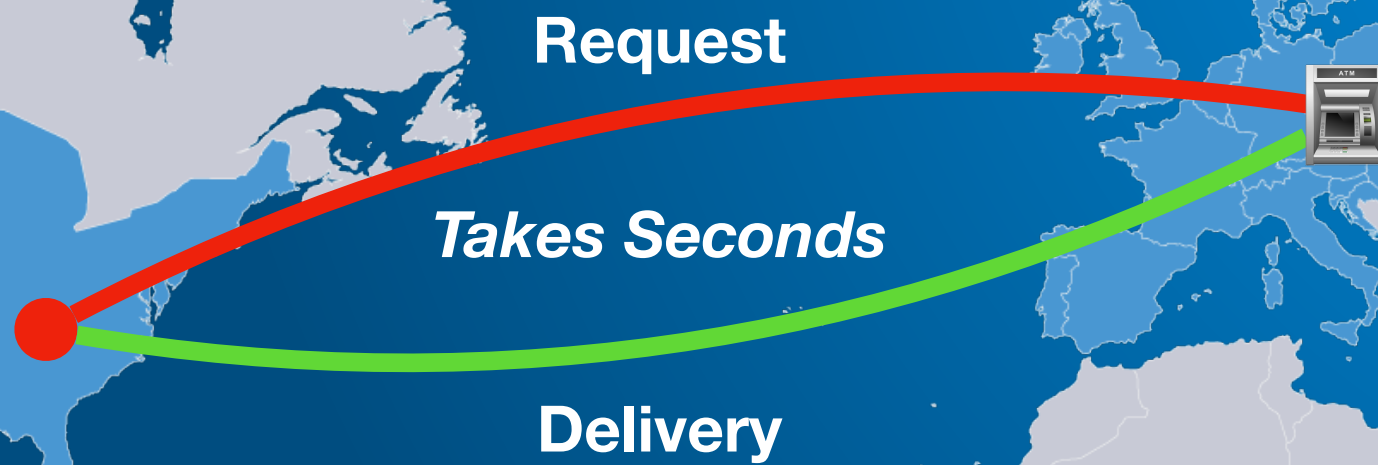


Delivery





# Analogy



# Analogy



Common stock issued by  
B&O Railroad Co. in 1903



Mortgage bond issued by Cleveland  
Short-Line Railway Co. in 1911



Promissory note issued by the  
Imperial Bank of India in 1926



Promissory note issued by the 2nd  
Bank of the United States in 1840



Common stock issued by  
Pennsylvania Railroad Co. in 1959



Common stock issued by  
Reading Co. in 1969



Mortgage note issued (signed) by  
"Shoeless" Joe Jackson in 1941



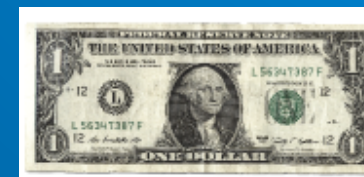
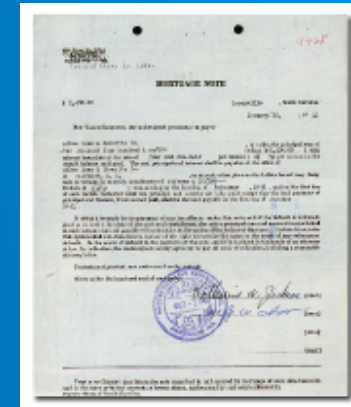
Certificate of deposit issued by the  
U.S. Postal Savings System in 1932



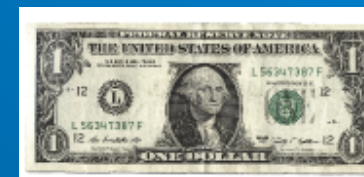
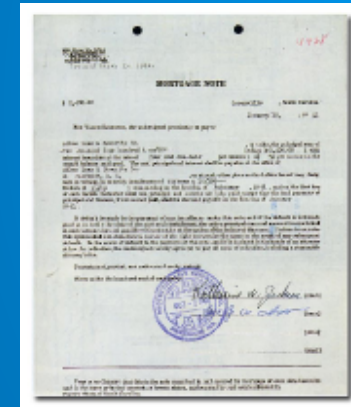
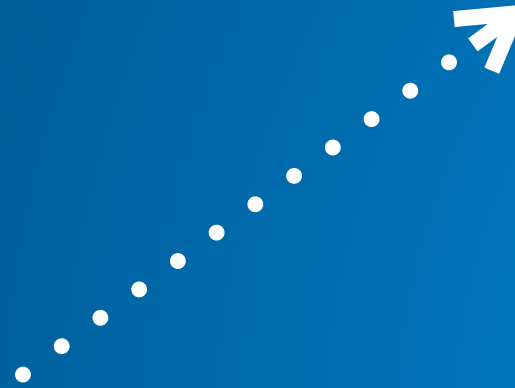
Federal reserve note issued by  
the U.S. Government in 2009



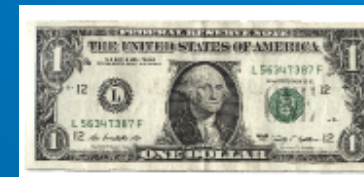
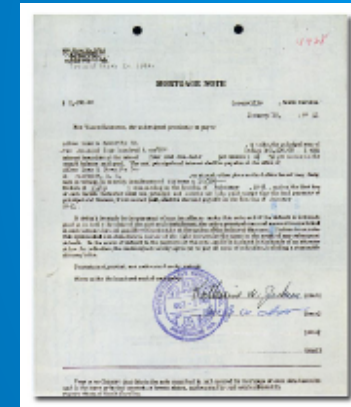
# Data Liquidity



# Data Liquidity

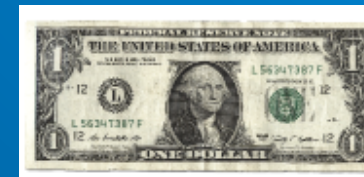
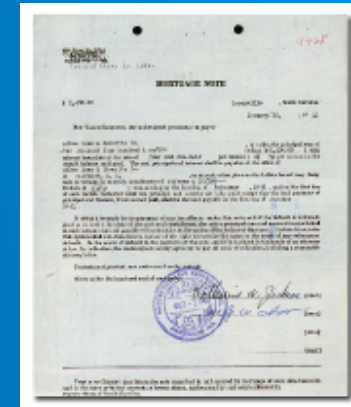


# Data Liquidity



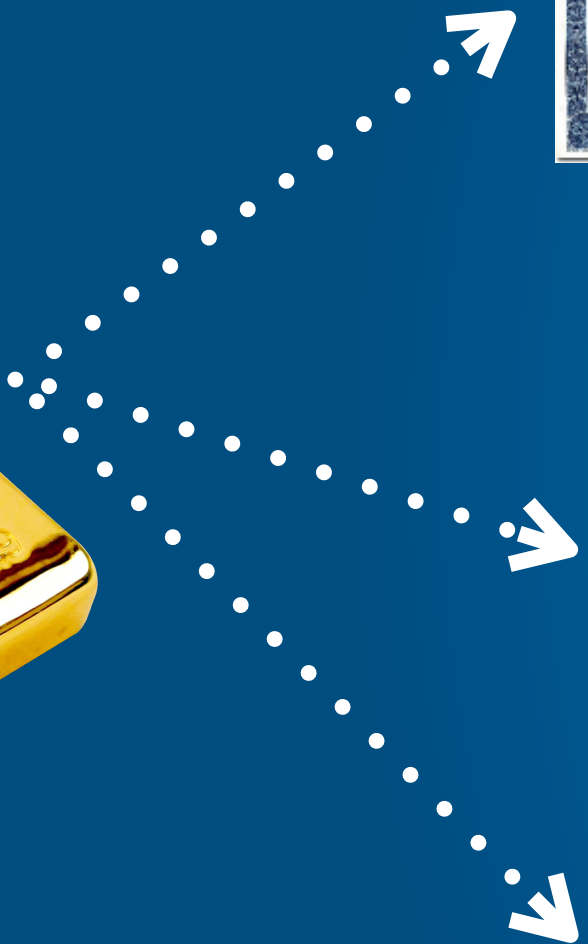


# Data Liquidity





# Data Liquidity



# Data Liquidity



# Data Liquidity

***" I don't know how many aortic stenosis patients I have "***

*Currently, it's difficult to even obtain basic counts of patients...*

# Data Liquidity

***" I don't know how many prostate cancer patients I have "***

*Currently, it's difficult to even obtain basic counts of patients...*



# Data Liquidity

*Beyond counts, what does data liquidity look like?*

# Data Liquidity

*Beyond counts, what does data liquidity look like?*

How many patients meet the eligibility criteria for this RSV study?

What is the door-to-balloon time for the past month?

How many MACE events did we see for this cohort in the last 6 months?



# Data Liquidity

*Beyond counts, what does data liquidity look like?*

How many patients meet the eligibility criteria for this RSV study?

What is the door-to-balloon time for the past month?

How many MACE events did we see for this cohort in the last 6 months?

I want to join liver MRIs with radiology and pathology reports for deep learning to predict hepatic cancer outcomes

# Data Liquidity

*Beyond counts, what does data liquidity look like?*

How many patients meet the eligibility criteria for this RSV study?

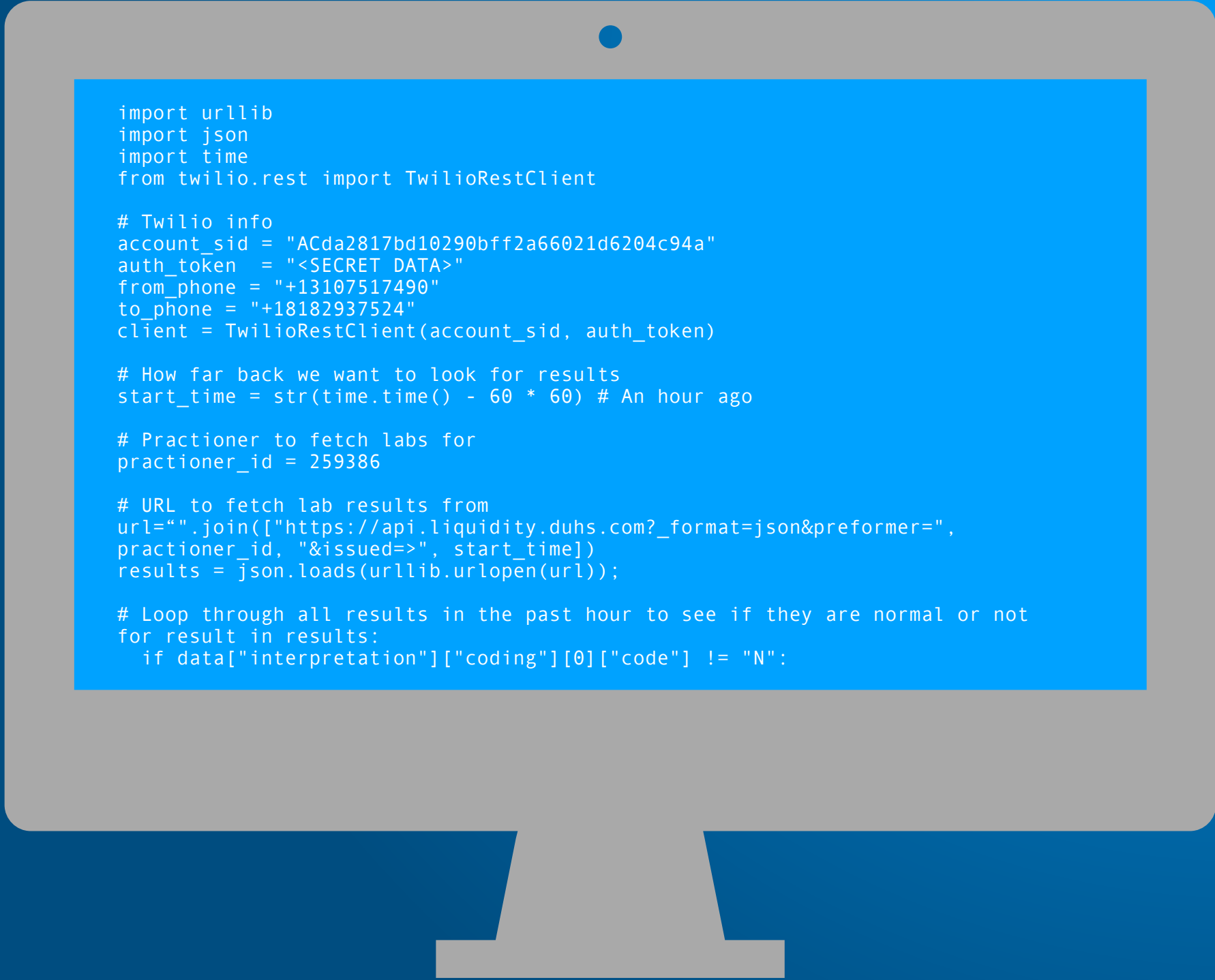
What is the door-to-balloon time for the past month?

How many MACE events did we see for this cohort in the last 6 months?

I want to join liver MRIs with radiology and pathology reports for deep learning to predict hepatic cancer outcomes

# Data Liquidity

*So what might liquidity look like?*



```
import urllib
import json
import time
from twilio.rest import TwilioRestClient

# Twilio info
account_sid = "ACda2817bd10290bff2a66021d6204c94a"
auth_token = "<SECRET DATA>"
from_phone = "+13107517490"
to_phone = "+18182937524"
client = TwilioRestClient(account_sid, auth_token)

# How far back we want to look for results
start_time = str(time.time() - 60 * 60) # An hour ago

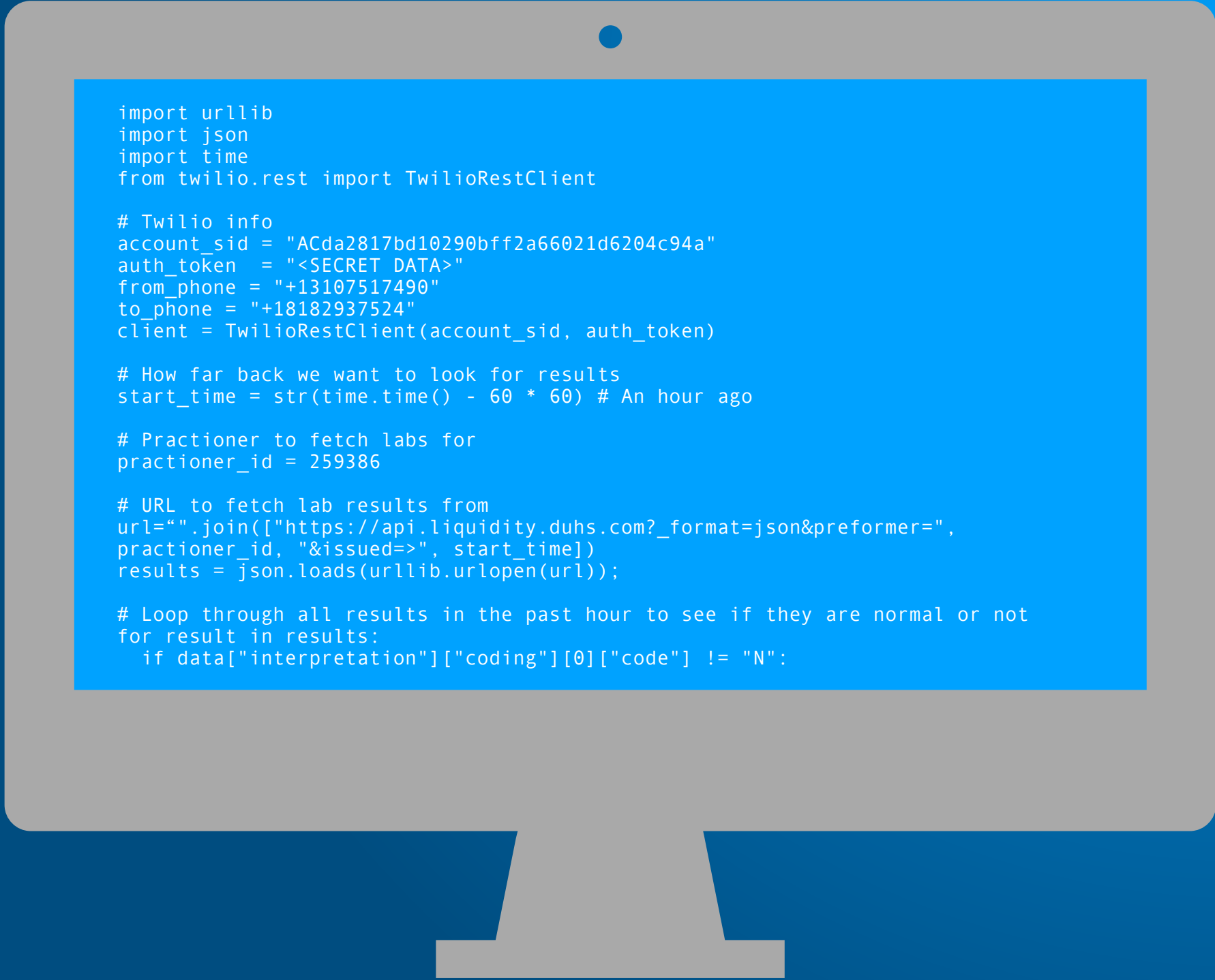
# Practioner to fetch labs for
practioner_id = 259386

# URL to fetch lab results from
url="".join(["https://api.liquidity.duhs.com?_format=json&preformer=",
practioner_id, "&issued=>", start_time])
results = json.loads(urllib.urlopen(url));

# Loop through all results in the past hour to see if they are normal or not
for result in results:
    if data["interpretation"]["coding"][0]["code"] != "N":
```

# Data Liquidity

*So what might liquidity look like?*



```
import urllib
import json
import time
from twilio.rest import TwilioRestClient

# Twilio info
account_sid = "ACda2817bd10290bff2a66021d6204c94a"
auth_token = "<SECRET DATA>"
from_phone = "+13107517490"
to_phone = "+18182937524"
client = TwilioRestClient(account_sid, auth_token)

# How far back we want to look for results
start_time = str(time.time() - 60 * 60) # An hour ago

# Practioner to fetch labs for
practioner_id = 259386

# URL to fetch lab results from
url="".join(["https://api.liquidity.duhs.com?_format=json&preformer=",
practioner_id, "&issued=>", start_time])
results = json.loads(urllib.urlopen(url));

# Loop through all results in the past hour to see if they are normal or not
for result in results:
    if data["interpretation"]["coding"][0]["code"] != "N":
```

# Data Liquidity

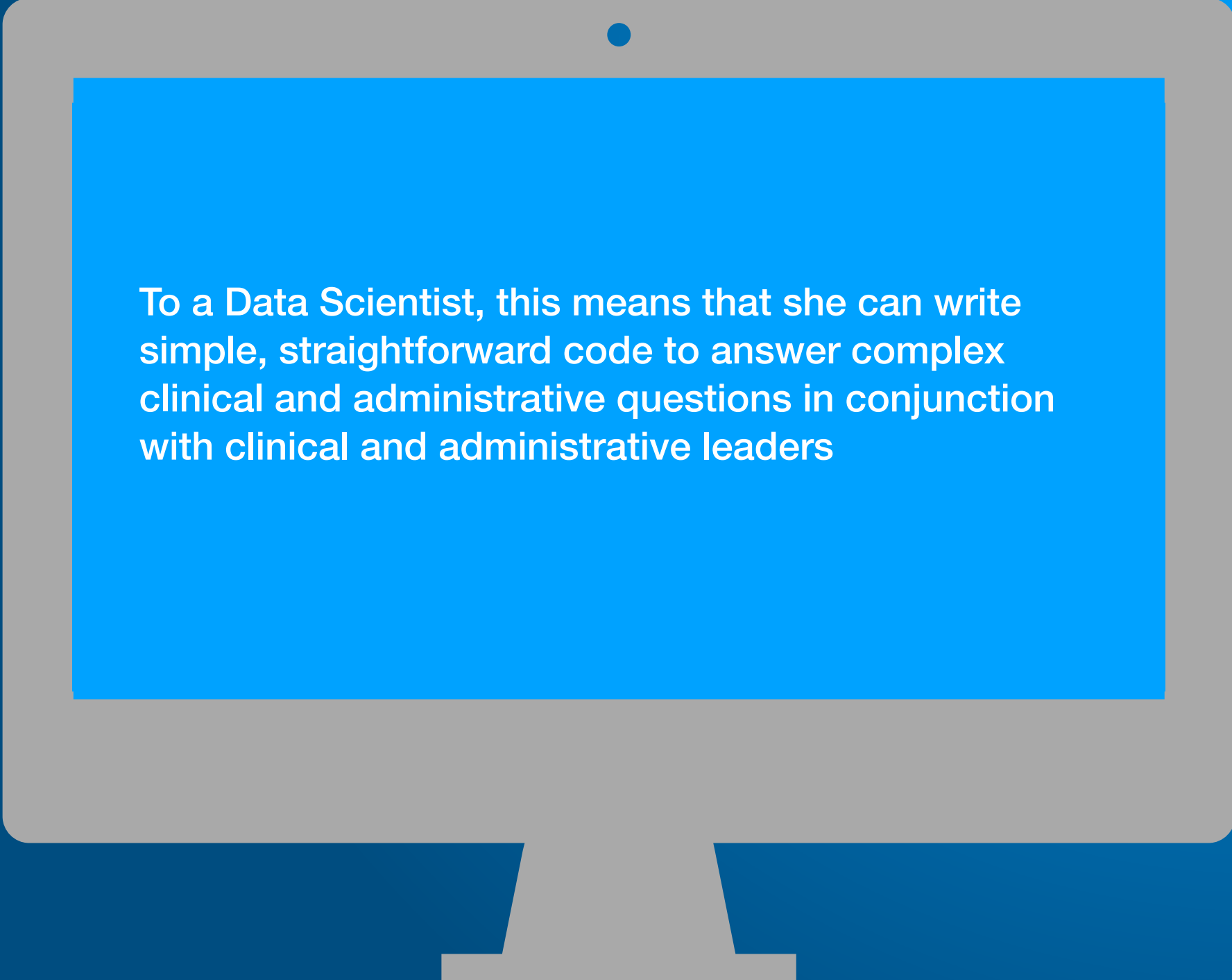
*So what might liquidity look like?*



To many, this Python code may look like gibberish.

# Data Liquidity

*So what might liquidity look like?*



To a Data Scientist, this means that she can write simple, straightforward code to answer complex clinical and administrative questions in conjunction with clinical and administrative leaders



# Data Liquidity

*To get there, we need to consider two components*

# Data Liquidity

*To get there, we need to consider two components*



# Policy

# Data Liquidity

*To get there, we need to consider two components*



## Policy



## Technology

# Data Liquidity

*before delving into those, it's probably helpful to discuss what “data liquidity” is **not***

Data Liquidity is not

*It is not untrammelled access to data  
for anyone regardless of its sensitivity*

Data Liquidity is not

*It is not untrammelled access to data  
for anyone regardless of its sensitivity*

*It is not solved only with technology*



# Data Liquidity is not

*It is not untrammelled access to data for anyone regardless of its sensitivity*

*It is not solved only with technology*

*It is not a pipe dream*

# Data Liquidity is

*Is agile & appropriate movement,  
merging, and analysis of data*

# Data Liquidity is

*Is agile & appropriate movement,  
merging, and analysis of data*

*Is where infrastructure & access are  
not the rate-limiting step*

# Data Liquidity is

***Is** agile & appropriate movement, merging, and analysis of data*

***Is** where infrastructure & access are not the rate-limiting step*

***Is** where analytic priorities, not process, drive use*

# Data Liquidity is

*Is agile & appropriate movement, merging, and analysis of data*

*Is where infrastructure & access are not the rate-limiting step*

*Is where analytic priorities, not process, drive use*

*Is secure, compliant, and auditable*



# Data Liquidity

*So let's consider these two components  
to data liquidity...*

# Data Liquidity

*So let's consider these two components  
to data liquidity...*



# Policy

# Data Liquidity

*So let's consider these two components to data liquidity...*



Policy

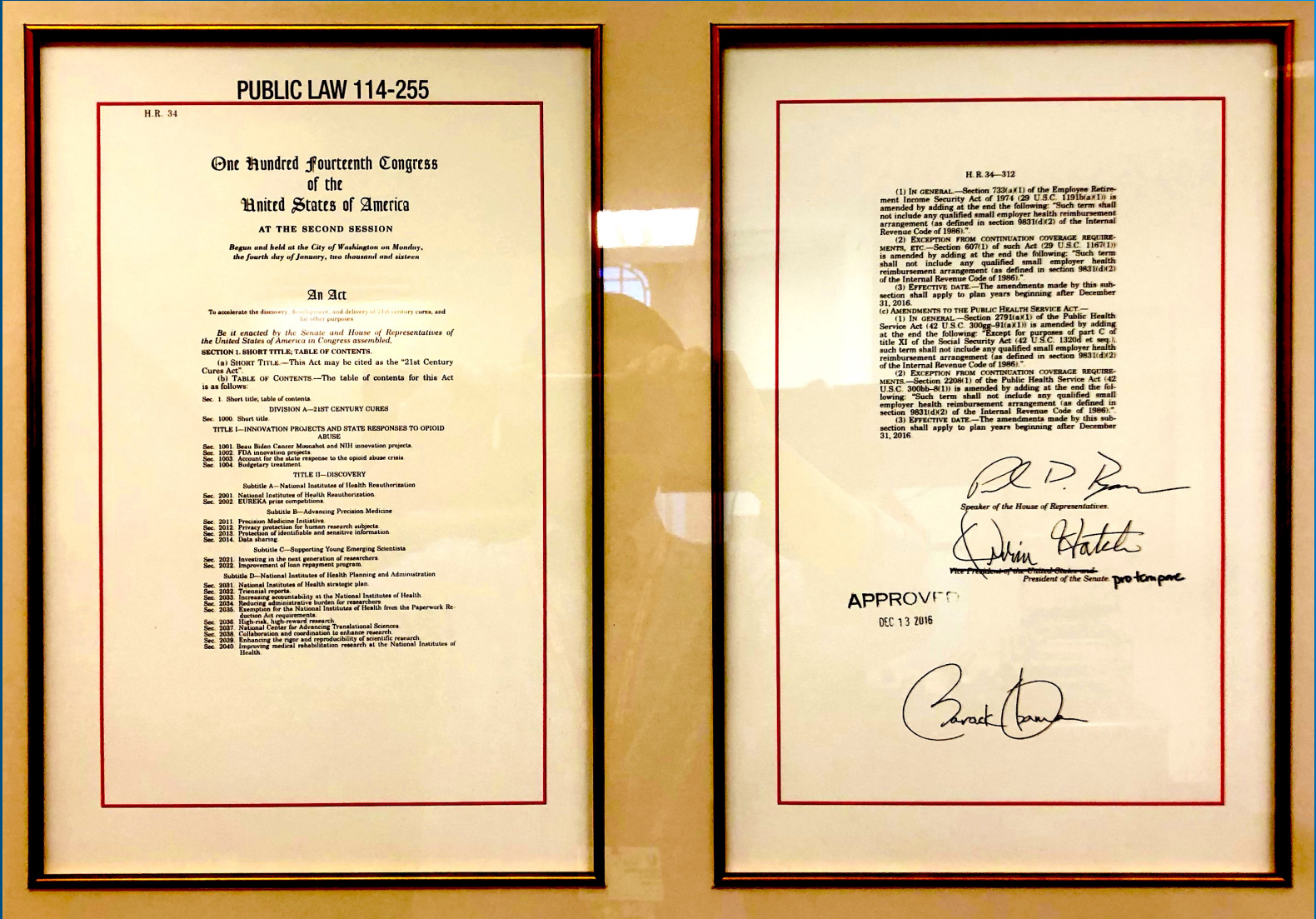


Technology

# Data Liquidity



# Policy





NATIONAL ARCHIVES

# FEDERAL REGISTER

The Daily Journal of the United States Government

[Sections](#)
[Browse](#)
[Search](#)
[Reader Aids](#)
[My FR](#)

0
[Sign in](#)
[Sign up](#)

---

**Rule**

---

## 21st Century Cures Act: Interoperability, Information Blocking, and the ONC Health IT Certification Program

---

A Rule by the [Health and Human Services Department](#) on [05/01/2020](#)

---

PUBLISHED DOCUMENT

☐ Start Printed Page 25642

**AGENCY:**

Office of the National Coordinator for Health Information Technology (ONC),  
Department of Health and Human Services (HHS).

**ACTION:**

Final rule.

**SUMMARY:**

This final rule implements certain provisions of the 21st Century Cures Act, including Conditions and Maintenance of Certification requirements for health information technology (health IT) developers under the ONC Health IT Certification Program (Program), the voluntary certification of health IT for use by pediatric health care providers, and reasonable and necessary activities that do not constitute information blocking. The implementation of these provisions will advance interoperability and support the access, exchange, and use of electronic health information. The rule also finalizes certain modifications to the 2015 Edition health IT certification criteria and Program in additional ways to

DOCUMENT DETAILS

**Printed version:**

[PDF](#)

**Publication Date:**

[05/01/2020](#)

**Agencies:**

[Department of Health and Human Services](#)  
[Office of the Secretary](#)

**Dates:**

Effective date: This final rule is effective on June 30, 2020.

**Effective Date:**

06/30/2020

**Document Type:**

Rule

**Document Citation:**

85 FR 25642

**Page:**

25642-25961 (320 pages)

**CFR:**

45 CFR 170



# Data Liquidity



# Policy

# Data Liquidity



## Policy

### 21st Century Cures Act, Section 4002

"... must also attest that it published application program interfaces ( ) and allows health information from such APIs to be accessible, exchanged and used without special effort through the use of APIs or successor technologies or standards, including providing access to all data elements of a patient's EHR to the extent permissible under applicable privacy laws."

# Data Liquidity

*Let's now look at*



# Technology

# Data Liquidity



# Technology

## *Data Standards*



*The diversity of data types: e.g. structured, unstructured, EHR, molecular, wearable, social determinants, &c...*



# Technology

## *Data Standards*



*The diversity of data types: e.g. structured, unstructured, EHR, molecular, wearable, social determinants, &c...*



*Formats not only have to account for structure, but transmissibility. Systems have to be prepared for realtime data transactions*





## Technology

### *Data Standards*



*The diversity of data types: e.g. structured, unstructured, EHR, molecular, wearable, social determinants, &c...*



*Formats not only have to account for structure, but transmissibility. Systems have to be prepared for realtime data transactions*



*Fungibility: machine learning learns with bulk data, but must be able to generate inference with individual data*

# Data Liquidity



# Technology

*Considerations include:*



*The diversity of data types: structured, unstructured, molecular, wearable, determinants, &c...*



*The velocity with which data originate and move. Systems prepared for realtime analysis*



*The scalability and flexibility of compute infrastructure to handle diverse data science workloads and these diverse data types*



# Data Liquidity



# Technology

*Considerations include:*



*Not a panacea!*



*The diversity of data types: structured, unstructured, molecular, wearable, determinants, &c...*



*The velocity with which data originate and move. Systems prepared for realtime analysis*



*The scalability and flexibility of compute infrastructure to handle diverse data science workloads and these diverse data types*

# Data Liquidity



# Technology

*Considerations include:*



- ★ *Immature ecosystem*
- ★ *More of a transaction standard*
- ★ *What's "inside the box" can be quite permissive*
- ★ *Needs real world critical mass for us to learn good vs bad implementations*

# Data Liquidity





# Data Liquidity



# Data Liquidity

*“the more that we use data, the clearer  
the river of data gets”*

*—Amy Abernethy*

# Future: Overcoming barriers to implementation

## Artificial Intelligence/Machine Learning (AI/ML)-Based Software as a Medical Device (SaMD) Action Plan

January 2021

Ref. 20

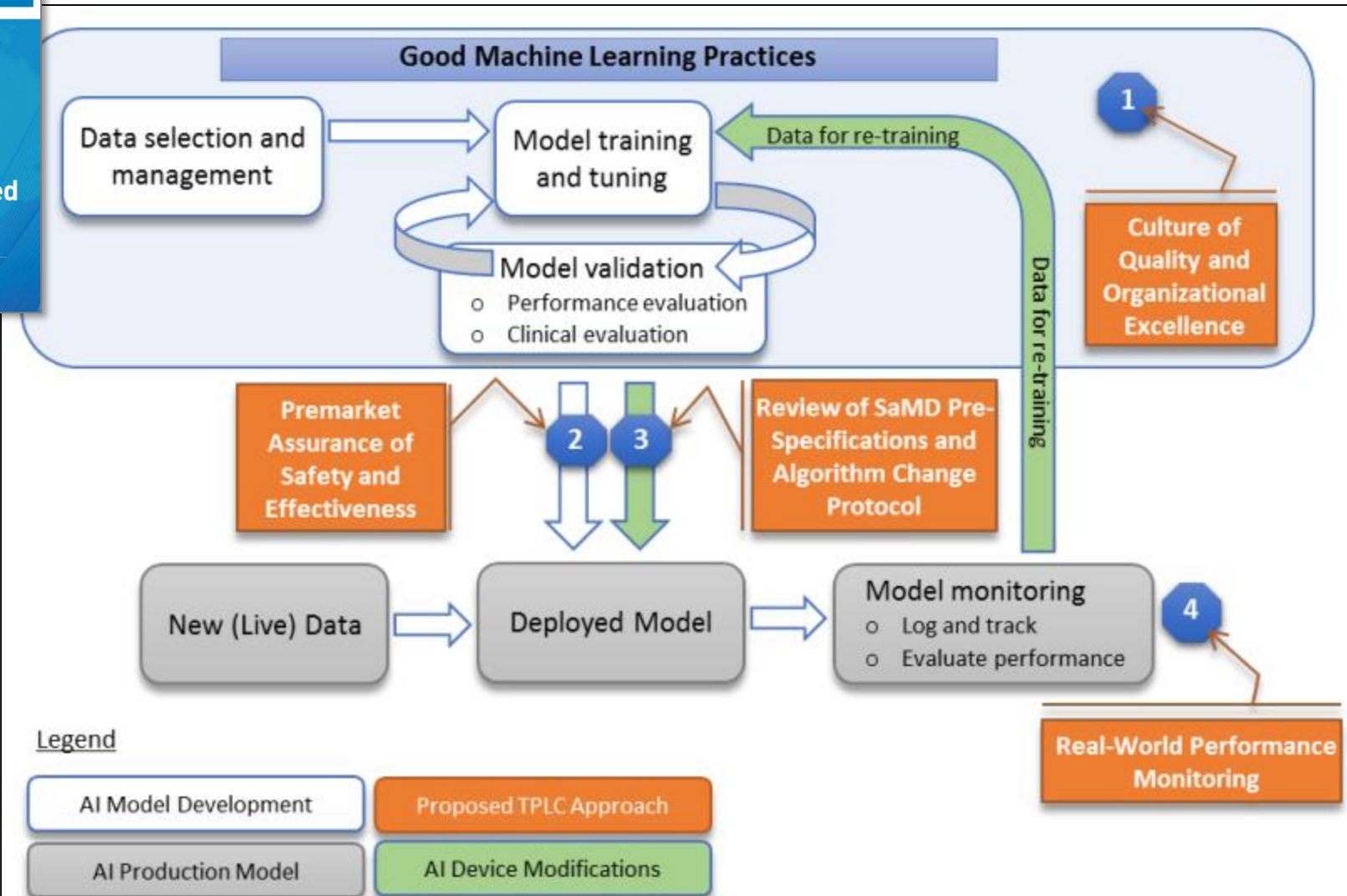
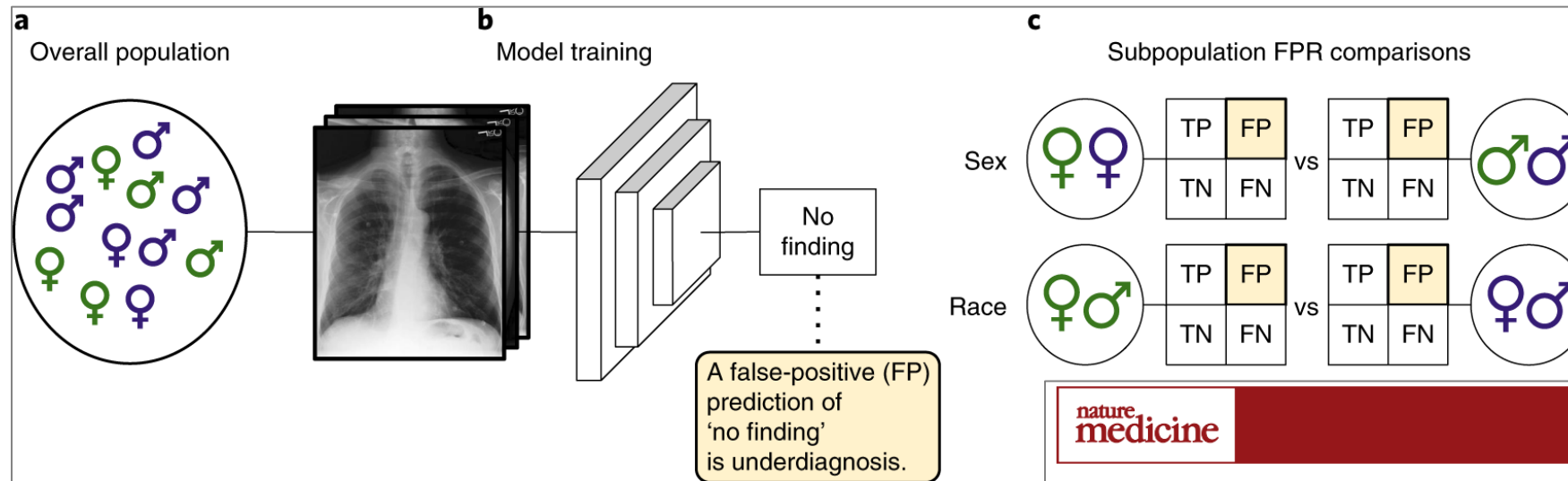


Figure 2: Overlay of FDA's TPLC approach on AI/ML workflow

# Future: Overcoming barriers to implementation



**nature**  
**medicine**

ARTICLES

<https://doi.org/10.1038/s41591-021-01595-0>

Check for updates

**OPEN**

**Underdiagnosis bias of artificial intelligence algorithms applied to chest radiographs in under-served patient populations**

Laleh Seyyed-Kalantari<sup>1,2</sup>✉, Haoran Zhang<sup>3</sup>, Matthew B. A. McDermott<sup>3</sup>, Irene Y. Chen<sup>3</sup> and Marzyeh Ghassemi<sup>2,3</sup>

# Future: Overcoming barriers to implementation

## Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings

Tolga Bolukbasi<sup>1</sup>, Kai-Wei Chang<sup>2</sup>, James Zou<sup>2</sup>, Venkatesh Saligrama<sup>1,2</sup>, Adam Kalai<sup>2</sup>

<sup>1</sup>Boston University, 8 Saint Mary's Street, Boston, MA

<sup>2</sup>Microsoft Research New England, 1 Memorial Drive, Cambridge, MA

tolgab@bu.edu, kw@kwchang.net, jamesyzou@gmail.com, srv@bu.edu, adam.kalai@microsoft.com

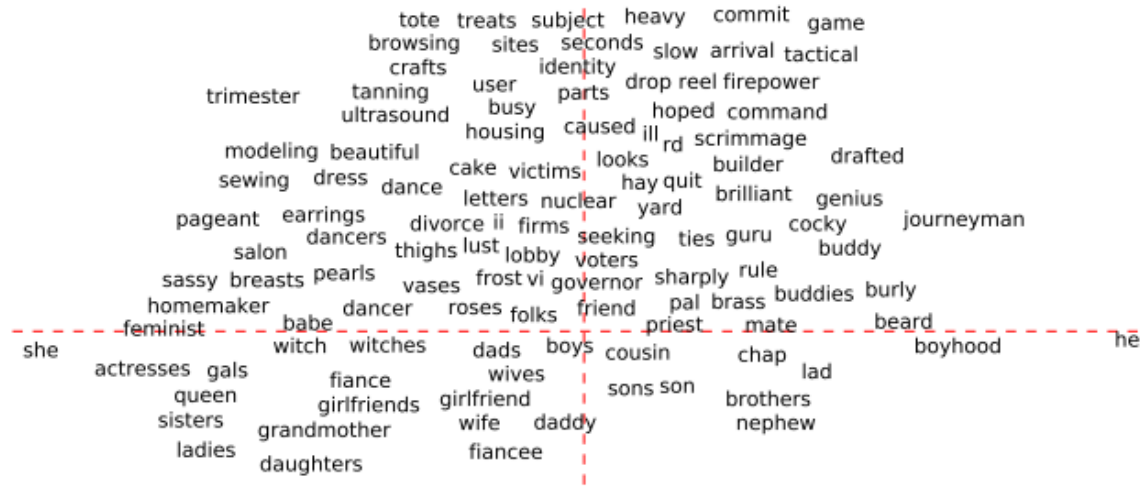


Figure 7: Selected words projected along two axes:  $x$  is a projection onto the difference between the embeddings of the words *he* and *she*, and  $y$  is a direction learned in the embedding that captures gender neutrality, with gender neutral words above the line and gender specific words below the line. Our hard debiasing algorithm removes the gender pair associations for gender neutral words. In this figure, the words above the horizontal line would all be collapsed to the vertical line.

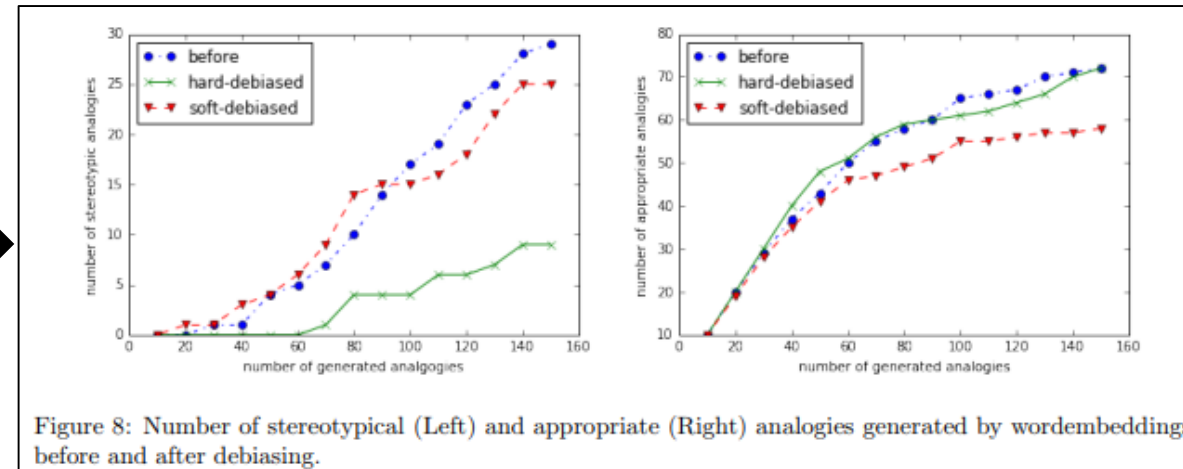
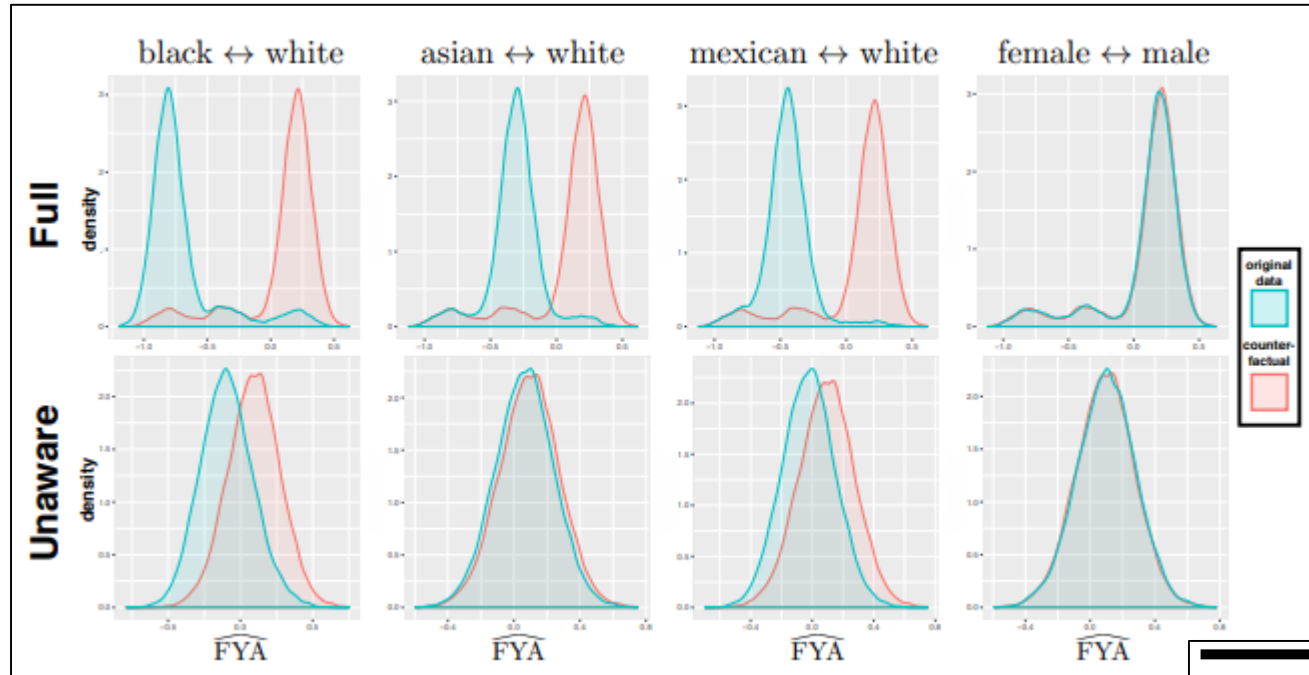


Figure 8: Number of stereotypical (Left) and appropriate (Right) analogies generated by word embeddings before and after debiasing.

Ref. 22



# Future: Overcoming barriers to implementation



## Counterfactual Fairness

**Matt Kusner \***  
The Alan Turing Institute and  
University of Warwick  
mkusner@turing.ac.uk

**Joshua Loftus \***  
New York University  
loftus@nyu.edu

**Chris Russell \***  
The Alan Turing Institute and  
University of Surrey  
crussell@turing.ac.uk

**Ricardo Silva**  
The Alan Turing Institute and  
University College London  
ricardo@stats.ucl.ac.uk

# Future: Overcoming barriers to implementation

**FAT / ML**

*fairness*  
*accountability*  
*transparency*

# **Future: Immediate next steps**

- **Educate yourself and trainees about ML techniques & reporting standards.**
- **Engage with efforts to define the regulatory perspective on ML in CR.**
- **Collaborate on proof-of-concept studies showing the promise of ML in CR & comparing ML to conventional approaches.**
- **Support data interoperability initiatives and advocate for patient-centered approaches to data ownership.**

# References

1. Weissler EH, Naumann T, Andersson T, et al. The role of machine learning in clinical research: transforming the future of evidence generation. *Trials*. Aug 16 2021;22(1):537. doi:10.1186/s13063-021-05489-x
2. Annapureddy AR, Angraal S, Caraballo C, et al. The National Institutes of Health funding for clinical research applying machine learning techniques in 2017. *NPJ Digit Med*. 2020;3:13. doi:10.1038/s41746-020-0223-9
3. Stevens LM, Mortazavi BJ, Deo RC, Curtis L, Kao DP. Recommendations for Reporting Machine Learning Analyses in Clinical Research. *Circ Cardiovasc Qual Outcomes*. Oct 2020;13(10):e006556. doi:10.1161/CIRCOUTCOMES.120.006556
4. Bisaso KR, Anguzu GT, Karungi SA, Kiragga A, Castelnuovo B. A survey of machine learning applications in HIV clinical research and care. *Comput Biol Med*. Dec 1 2017;91:366-371. doi:10.1016/j.combiomed.2017.11.001
5. Kim KJ, Tagkopoulos I. Application of machine learning in rheumatic disease research. *Korean J Intern Med*. Jul 2019;34(4):708-722. doi:10.3904/kjim.2018.349
6. Wong CH, Siah KW, Lo AW. Estimation of clinical trial success rates and related parameters. *Biostatistics*. Apr 1 2019;20(2):273-286. doi:10.1093/biostatistics/kxx069
7. Fauqueur J TA, Togia T. Constructing large scale biomedical knowledge bases from scratch with rapid annotation of interpretable patterns. *Proceedings of the 18th BioNLP Workshop and Shared Task*. 2019;doi:doi:10.18653/v1/w19-5016.
8. Madhukar NS, Khade PK, Huang L, et al. A Bayesian machine learning approach for drug target identification using diverse data types. *Nat Commun*. Nov 19 2019;10(1):5221. doi:10.1038/s41467-019-12928-6
9. Liu Q AM, Brockschmidt M, Gaunt AL. Constrained graph variational autoencoders for molecule design. *NeurIPS 2018*. 2018;doi:arXiv:1805.09076
10. Langner S, Hase F, Perea JD, et al. Beyond Ternary OPV: High-Throughput Experimentation and Self-Driving Laboratories Optimize Multicomponent Systems. *Adv Mater*. Apr 2020;32(14):e1907801. doi:10.1002/adma.201907801
11. Zhao Y, Zeng D, Socinski MA, Kosorok MR. Reinforcement learning strategies for clinical trials in nonsmall cell lung cancer. *Biometrics*. Dec 2011;67(4):1422-33. doi:10.1111/j.1541-0420.2011.01572.x
12. Seymour CW, Kennedy JN, Wang S, et al. Derivation, Validation, and Potential Treatment Implications of Novel Clinical Phenotypes for Sepsis. *JAMA*. May 28 2019;321(20):2003-2017. doi:10.1001/jama.2019.5791
13. Glicksberg BS, Miotto R, Johnson KW, et al. Automated disease cohort selection using word embeddings from Electronic Health Records. *Pac Symp Biocomput*. 2018;23:145-156.
14. Chen R, Jankovic F, Marinsek N, et al. Developing Measures of Cognitive Impairment in the Real World from Consumer-Grade Multimodal Sensor Streams. presented at: Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery; Data Mining; 2019; Anchorage, AK, USA. <https://doi.org/10.1145/3292500.3330690>
15. Loudon K, Treweek S, Sullivan F, Donnan P, Thorpe KE, Zwarenstein M. The PRECIS-2 tool: designing trials that are fit for purpose. *BMJ*. May 8 2015;350:h2147. doi:10.1136/bmj.h2147
16. Schork NJ. Personalized medicine: Time for one-person trials. *Nature*. Apr 30 2015;520(7549):609-11. doi:10.1038/520609a
17. Mitchell M, Wu S, Zaldivar A, et al. Model Cards for Model Reporting. presented at: Proceedings of the Conference on Fairness, Accountability, and Transparency; 2019; Atlanta, GA, USA. <https://doi.org/10.1145/3287560.3287596>
18. Sterckx S, Rakic V, Cockbain J, Borry P. "You hoped we would sleep walk into accepting the collection of our data": controversies surrounding the UK care.data scheme and their wider relevance for biomedical research. *Med Health Care Philos*. Jun 2016;19(2):177-90. doi:10.1007/s11019-015-9661-6
19. Robbins RB, E. An invisible hand: Patients aren't being told about the AI systems advising their care. *STAT*. 15 July 2020. <https://www.statnews.com/2020/07/15/artificial-intelligence-patient-consent-hospitals/>
20. Food and Drug Administration. Artificial Intelligence/Machine Learning (AI/ML)-Based Software as a Medical Device (SaMD) Action Plan (2021). <https://www.fda.gov/media/145022/download>
21. Seyyed-Kalantari L, Zhang H, McDermott MBA, Chen IY, Ghassemi M. Underdiagnosis bias of artificial intelligence algorithms applied to chest radiographs in under-served patient populations. *Nat Med*. Dec 2021;27(12):2176-2182. doi:10.1038/s41591-021-01595-0
22. Bolukbasi T, Chang K-W, Zou J, Saligrama V, Kalai A. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. presented at: Proceedings of the 30th International Conference on Neural Information Processing Systems; 2016; Barcelona, Spain.
23. Kusner MJ, Russell C, Silva R. Counterfactual fairness. *Adv Neural Inf Process Syst*. 2017;doi:arXiv:1703.06856



# Thank you

Questions and comments?