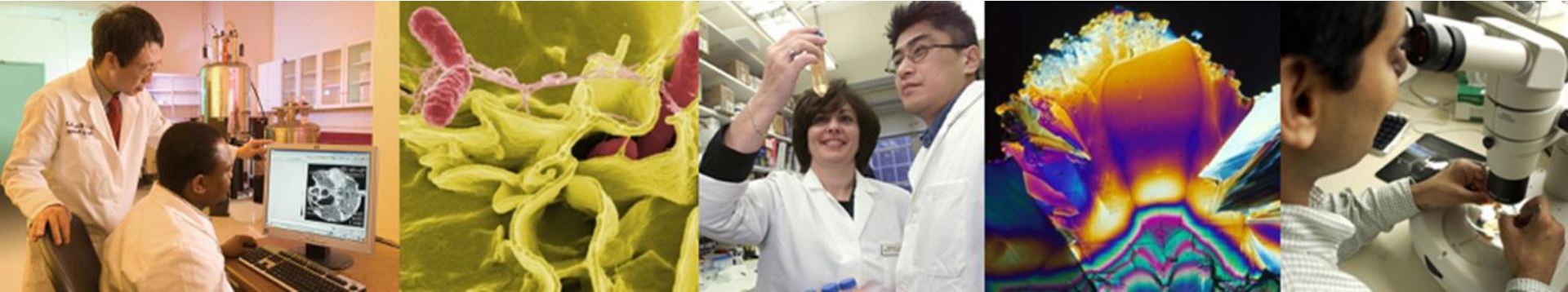# Artificial Intelligence Working Group Update

118th Meeting of the Advisory Committee to the Director (ACD)
*June 13, 2019*



**David Glazer**
Engineering Director, Verily

**Lawrence A. Tabak, DDS, PhD**
Principal Deputy Director, NIH
Department of Health and Human Services

# Agenda

- **refresher on [our charge](#)**
- interim progress report
- next steps

# We Generate Enormous Volumes of Data Daily



FROM HEAD TO TOE WEARABLE TECHNOLOGY

**SHIRT**
Conductive thread means a computer is literally built into the fabric of the shirt, providing the processing power for all the other wearable gadgets.

**WRISTBAND**
A sensor that tracks movement to determine the number of steps taken through the day – 10,000 is ideal – and how much sleep the wearer gets at night.

**TROUSERS**
Also made with conductive thread, the trousers take the energy generated by movement and use it to power the other gadgets.

**GLASSES**
Overlays navigation directions and information about points of interest directly on to the wearer's field of vision.

**WRISTWATCH**
Vibrates when a message arrives and displays it on the watch face. Tells the time too.

**HAND**
Embedded under the skin is a chip containing medical records, passport data and credit records. Information is transferred by waving the hand over a suitable scanner.

**SHOES**
GPS chip provides directions using LED lights in each shoe: the left shoe indicates direction, while the right shows distance.

GRAPHIC: JOHN BRADLEY

https://people.rit.edu/sml2565/iimproject/wearables/index.html

3

# Data Science at NIH: A Snapshot

- CIT supports a 100GB Network moving 4PB of data per day
- Datasets and resources
  - List of extramural programs generating datasets (only a subset)
  - Datasets supported across IC and topic area
  - Range in size from several hundred terabytes to several petabytes
    - SRA and dbGaP, ~15 PB of genomic sequence data
      - Controlled access ~8 PB
      - Open access ~6 PB
    - GTEx, ~200 TB

| DATASET | Primary IC |
| --- | --- |
| ABCD (Adolescent Brain Cognitive Development) | MH |
| Accelerating Medicine Partnership - Parkinson's Disease (AMP PD) | NS |
| Age-Related Eye Disease Study (AREDS2) | EY |
| All of Us Research Program | OD |
| BRAIN Initiative | many |
| Biomedical Translational Research Information System (BTRIS) | CC |
| dbGAP | NL |
| Framingham Studies | HL |
| Gabriella Miller Kids First Pediatric Research Program | CF/HL |
| Genotype-Tissue Expression (GTEx) | CF/HG |
| Cancer Genome Characterization Initiative (CGCI) | NCI |
| Analysis, Visualization, and Informatics Lab-space (AnVIL) | HG |
| Chest and Cardiac Image Archive | HL |
| Genetics of Alzheimer's Disease Project (NIAGADS) | NIA |
| RSNA Radiology Image Share | EB |
| The Cancer Genome Atlas Project (TCGA) | NCI |
| TOPMed | HL |
| Alliance for Genome Resources Model Organism Databases (MODs) | HG |
| ClinVar | NL |
| dbSNP | NL |
| ENCODE | HG |
| Gene Expression Omnibus (GEO) | NL |
| MACS/WIHS Longitudinal AIDS Data | AI |
| Neuroimaging Tools & Resources (NITRC) | EB |
| SRA | NL |
| UniProt | HG/GM |

# Every Day Artificial Intelligence Applications

**ARTIFICIAL INTELLIGENCE**
A program that can sense, reason, act, and adapt
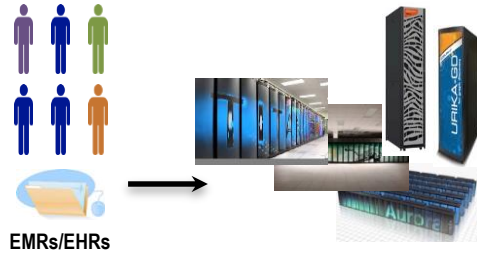
**MACHINE LEARNING**
Algorithms whose performance improve as they are exposed to more data over time
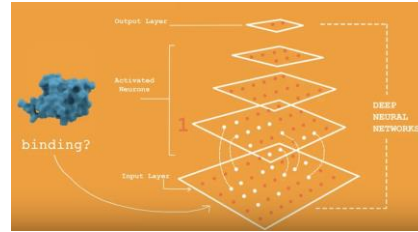
**DEEP LEARNING**
Subset of machine learning in which multilayered neural networks learn from vast amounts of data

# AI in Biomedicine: Opportunities



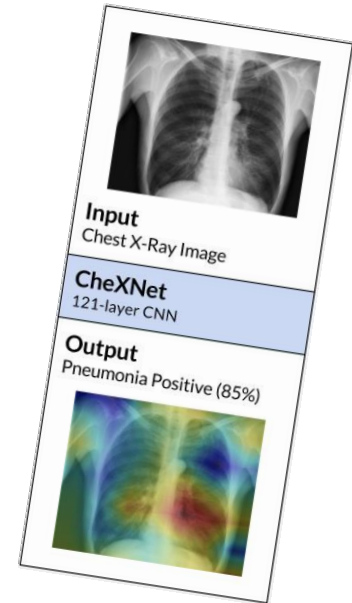Extract medical information from text in EMRs/EHRs

Interpret genomic sequence data to understand impact of mutations on protein function

Read medical images and help diagnose diseases like pneumonia and cancer

**Input**
Chest X-Ray Image

**CheXNet**
121-layer CNN

**Output**
Pneumonia Positive (85%)

Monitor sleep and vitals to send information about health at home to doctors

Determine which calls to child welfare systems warrant deployment of family support and prevention resources to protect at-risk children
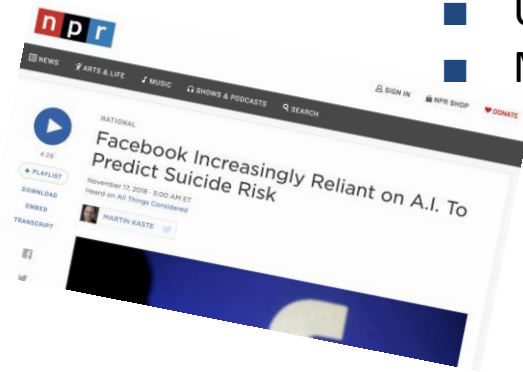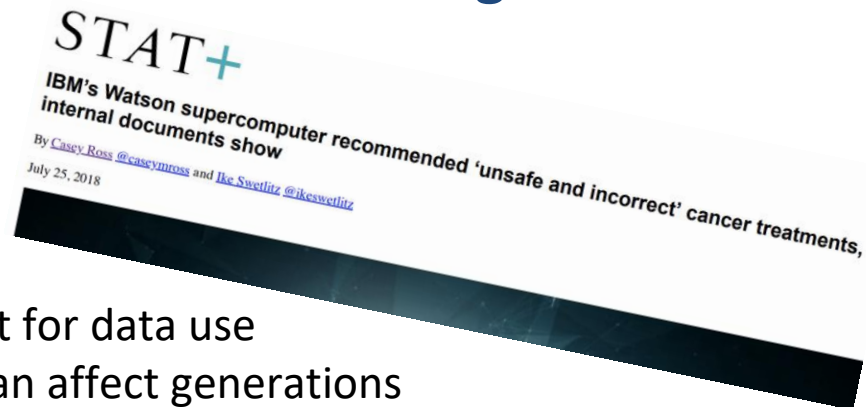
7

Examples from Katabi, Ng, Putnam-Hornstein, Troyanskaya, and others

# AI in Biomedicine: Legal and Ethical Challenges



PROPUBLICA

**Sloan Kettering's Cozy Deal With Start-Up Ignites a New Uproar**

A for-profit venture with exclusive rights to use the cancer center's vast archive of tissue slides has generated concerns among pathologists at the hospital, as well as experts in nonprofit law and corporate governance.

by Charles Ornstein, ProPublica, and Katie Thomas, The New York Times, Sept. 20, 2018, 4:10 p.m. EDT



*STAT+*

**IBM's Watson supercomputer recommended 'unsafe and incorrect' cancer treatments, internal documents show**

By Casey Ross @caseymross and Ike Swetlitz @ikeswetlitz

July 25, 2018

- No clear rules on consent for data use
- Threats to privacy that can affect generations
- How can people opt out? at the beginning or later on?
- Potential for bias and discrimination
- Use of incomplete or selective data
- Misuse of data



NATIONAL

**Facebook Increasingly Reliant on A.I. To Predict Suicide Risk**

November 17, 2018 · 5:00 AM ET
Heard on All Things Considered

MARTIN KASTE



Opinion | THE PRIVACY PROJECT    The New York Times

**Insurers Want to Know How Many Steps You Took Today**

The cutting edge of the insurance industry involves adjusting premiums and policies based on new forms of surveillance.

By Sarah Jeong

Ms. Jeong is a member of the editorial board.

**Controversy at MSK Cancer Center Regarding the Pathology Archive and Database**

# Charge to the AI Working Group (December 14, 2018)

- Are there opportunities for cross-NIH effort in AI? How could these efforts reach broadly across biomedical topics and have positive effects across many diverse fields?

- How can NIH help build a bridge between the computer science community and the biomedical community?

- What can NIH do to facilitate training that marries biomedical research with computer science?
  - Computational and biomedical expertise are both necessary, but careers may not look like traditional tenure track positions that follow the path from PhD to post-doc to faculty

- Identify the major ethical considerations as they relate to biomedical research and using AI/ML/DL for health-related research and care, and suggest ways that NIH can build these considerations into its AI-related programs and activities

# ACD Artificial Intelligence Working Group Members

Rediet Abebe
Cornell

Greg Corrado,
PhD
Google

Kate Crawford,
PhD
AI Now Institute

Barbara
Engelhardt, PhD
Princeton

David Glazer
Verily (Co-Chair)

David
Haussler, PhD
USCS

Dina Katabi, PhD
MIT Computer
Science & AI Lab

Daphne
Koller, PhD
insitro

Anshul Kundaje,
PhD
Stanford University

Eric Lander, PhD
Broad Institute

Jennifer
Listgarten, PhD
Berkeley

Michael
McManus, PhD
Intel

Lawrence Tabak,
DDS, PhD
NIH (Co-Chair)

Serena
Yeung, PhD
Harvard

# Agenda

- refresher on our charge
- **interim progress report**
- next steps

# Themes

1. more AI-ready **data**
2. more **multilingual** researchers
3. **ELSI**: ethical, legal, and social implications
4. important areas to **apply** AI
5. important areas to **advance** AI

For each theme:

- Opportunity -- why this area is important
- Do Now -- recommendations for action in H2 2019
- Questions --  what we'll be drilling into next

| data | multilingual | ELSI | apply | advance |

# more AI-ready data

Opportunity

- Creation and stewardship of data sets to enable machine learning may be the NIH's single greatest lever to accelerate the application of AI within biomedicine.
- Such datasets must be deliberately constructed to be useful, representative, and ethical. The responsible distribution and maintenance of these data sets is as important as their creation.

*"Show me the data."*
*- every ML researcher ever*

| data | multilingual | ELSI | apply | advance |

# more AI-ready <u>data</u>

<u>Do Now</u>

- Begin cataloging existing datasets, tracking attributes relevant for AI-readiness.

<u>Questions to answer next</u>

- What are the attributes of AI-ready data?
  - e.g. rich -- multi-modal, longitudinal, well-labeled
  - e.g. usable -- friendly data access, harmonized with other datasets
  - e.g. beyond observation -- includes perturbation data
- How can NIH accelerate existence of and promote access to AI-ready data?
- How can NIH nurture novel sources of health data?
  - e.g. by expanding the use of sensors to create training datasets
  - e.g. by building libraries of genomic and other molecular data
  - e.g. by enabling sharing of data derived from care delivery
  - e.g. by applying aggressive identity-protecting techniques to real data

| data | multilingual | ELSI | apply | advance |
|------|--------------|------|-------|---------|

# more <u>multilingual</u> researchers

<u>Opportunity</u>

- Increase the pool of people who are "bilingual," meaning they have experience in both biomedicine and computer science, expertise in at least one, and can bridge the two worlds.

- Broaden the tent further, to rigorously address representation and equity. Include social and behavioral scientists who study health disparities and other issues in populations that are underrepresented in biomedical research.

*More perspectives ⇒ better results*

| data | multilingual | ELSI | apply | advance |

# more <u>multilingual</u> researchers

<u>Do Now</u>

- Co-sponsor workshop proposal on *Learning Meaningful Representations of Life* at <u>NeurIPS</u> (leading ML conference) in December 2019
    - goal is deeper NIH connections to ML community; perhaps keynote by Francis Collins?
- Allocate at least one third of this year's fellows slots to AI projects
    - in the NIH Civic Digital Fellows and the OD's National Service Sabbatical data program

<u>Questions to answer next</u>

- What are the attributes of multilingual researchers?
- How can NIH upgrade curricula to train multilingual researchers?
    - all levels of education -- secondary, undergrad, graduate, professional
- How can NIH catalyze an active community of multilingual researchers?

| data | multilingual | ELSI | apply | advance |

# <u>ELSI</u>: ethical, legal, and social implications

<u>Opportunity</u>

- This is one of the biggest challenge areas for biomedical applications of AI, since inappropriate use can present real harms, especially to under-represented and marginalized populations.

- Much more work is needed on building the guardrails to ensure safety, ethical deployment, and non-discriminatory impacts. NIH can set the quality standard, develop more rigorous frameworks around potential harms and challenges, and create the world's best safeguards.

- Above all, NIH can take a leadership position in building strong oversight and accountability mechanisms for the use of AI in biomedicine.

*These tools have sharp edges -- let's "do no harm".*

| data | multilingual | ELSI | apply | advance |

# ELSI: ethical, legal, and social implications

<u>Do Now</u>

- Draft a charter for an NIH governance / advisory body on AI standards and ethics

<u>Questions to answer next</u>

- How can NIH set standards for labeling training data? (see <u>datasheets</u> paper)
    - like Rx labels: "here's what you should know before taking this data"
    - e.g. data sourcing, relevant ethical/legal topics, (in)appropriate use
- How can NIH set standards for labeling ML models? (see <u>model cards</u> paper)
    - e.g. where should / shouldn't it be used?
    - a model's label should include the labels of its training data
- How can NIH update/improve ethical review processes to guide suitable use of AI in biomedical research?
- How can NIH help educate the community on the social and legal risks of AI?

| data | multilingual | ELSI | apply | advance |

# important areas to <u>apply</u> AI

<u>Opportunity</u>

- There are numerous directions in biomedical research, public health, and healthcare management where advances in AI are underutilized, yet if they are integrated with these fields, they could potentially lead to transformative impact.

- NIH can encourage the exploration of the above AI applications, and facilitate interactions between AI experts and researchers in the fields of biomedical research, public health, and healthcare delivery and management.

*The tech community doesn't know where to help -- let's tell them.*

| data | multilingual | ELSI | apply | advance |

# important areas to <u>apply</u> AI

<u>Do Now</u>

- Create a "top 10" catalog of success stories in this space, and use to inspire future investment and creativity.
    - e.g. recent [FDA AI approvals](#)

<u>Questions to answer next</u>

- How do we recognize the biomedical and public health opportunities that would benefit most from the application of AI?
    - E.g. specific disease or diagnostic areas
    - E.g. reduce health disparities
- How do we raise awareness of those opportunities, and catalyze problem-solving collaborations between MDs and AI experts?
    - E.g. workshops on specific subareas, such as impact on minority populations
    - E.g. well-defined problems and ways to measure success vs. state of the art

| data | multilingual | ELSI | apply | advance |

# important areas to <u>advance</u> AI

<u>Opportunity</u>

- Realizing the full potential of AI in biology, medicine and healthcare requires advancing AI beyond current capabilities, and solving some of the biggest challenges and open problems in AI. Biomedicine and healthcare are therefore valuable domains for motivating and grounding fundamental research in AI.

*New hard problems need new powerful tools.*

| data | multilingual | ELSI | apply | advance |
|------|--------------|------|-------|---------|

# important areas to <u>advance</u> AI

<u>Do Now</u>

- n/a

<u>Questions to answer next</u>

- How can NIH support the development of methods that:
  - can learn effectively in unlabeled, weakly labeled, and semi-supervised regimes
  - can perform challenging tasks beyond pattern recognition and supervised learning
  - can generalize predictions to real-world scenarios not reflected in training datasets
- How can NIH support the development of algorithms that:
  - are interpretable and safe
  - can effectively integrate with large, real-world systems

| data | multilingual | ELSI | apply | advance |
|------|--------------|------|-------|---------|

# Agenda

- refresher on our charge
- interim progress report
- **next steps**

# "Do Now" recommendations

| | |
|---|---|
| AI-ready data | ▪ Begin cataloging existing datasets, tracking attributes relevant for AI-readiness. |
| multilingual researchers | ▪ Co-sponsor workshop proposal on *Learning Meaningful Representations of Life* at NeurIPS.<br>▪ Allocate at least one third of this year's fellows slots to AI projects. |
| ELSI | ▪ Draft a charter for an NIH governance / advisory body on AI standards and ethics. |
| areas to apply AI | ▪ Create a "top 10" catalog of success stories in this space, and use to inspire future investment and creativity. |

# Questions to answer next

| | |
|---|---|
| AI-ready data | ▪ What are the attributes of AI-ready data?<br>▪ How can NIH accelerate existence of and promote access to AI-ready data?<br>▪ How can NIH nurture novel sources of health data? |
| multilingual researchers | ▪ What are the attributes of multilingual researchers?<br>▪ How can NIH upgrade curricula to train multilingual researchers?<br>▪ How can NIH catalyze an active community of multilingual researchers? |
| ELSI | ▪ How can NIH set standards for labeling training data?<br>▪ How can NIH set standards for labeling ML models?<br>▪ How can NIH update/improve ethical review processes to guide suitable use of AI?<br>▪ How can NIH help educate the community on the social and legal risks of AI? |
| areas to apply AI | ▪ How do we recognize the biomedical and public health opportunities that would benefit most from the application of AI?<br>▪ How do we raise awareness of those opportunities, and catalyze problem-solving collaborations between MDs and AI experts? |
| areas to advance AI | ▪ How can NIH support the development of novel methods?<br>▪ How can NIH support the development of novel algorithms? |

# Timeline

- Feb 2019: kickoff meeting
- **June 2019: Interim report to ACD**
- Dec 2019: final recommendations to ACD
- beyond 2019: the group will convene intermittently, as needed but infrequently, for updates and continued guidance



Serena Yeung (Stanford)
Michael (Intel)
Eric Lander (Broad)
Anshul Kundaje (Stanford)
Daphne Koller (insitro)
Barbara Engelhardt (Princeton)
Kate Crawford (NYU)
Rediet Abebe (Cornell)
Greg Corrado (Verily)
David Glazer (Verily)
Francis Collins (NIH)
Dina Katabi (MIT)
Jessica Mazerik (NIH)
Jennifer Listgarten, (Berkeley)

AI WG kick-off meeting at NIH in February 2019

# NIH…
## Turning Discovery Into Health

Lawrence.Tabak@nih.gov