

MACHINE LEARNING IN THE OPERATING ROOM

FRANK RUDZICZ



UNIVERSITY OF
TORONTO

St. Michael's

Inspired Care.
Inspiring Science.



**SURGICAL SAFETY
TECHNOLOGIES**



**Standards Council of Canada
Conseil canadien des normes**

TODAY

1. Some risks of deep machine learning
 - a) Mitigation of those risks with 'Explainable AI'
2. Potential of deep machine learning in the operating room
3. Implications for practice

STOKING THE HYPE: TO ERR IS HUMAN?

- **Humans** are notoriously **bad** with information.
 - Patients **misread** or **miscommunicate** their own symptoms.
 - Nearly **half** of American adults have difficulty understanding and acting upon health information (IOM, 2004).
 - Faulty memory; skill obsolescence; cognitive biases; cognitive/time limitations; **recency biases**; other human biases.
 - *Diagnoses* correlate with advertising and media exposure.
- Winters *et al.* (2012) showed that ~40,500 patients die in ICU, in the USA, each year due to misdiagnosis.

STOKING THE HYPE: TO ERR IS HUMAN?

- Graber et al. (2005) studied one hundred cases of **diagnostic error** involving internists ...
 - **Cognitive factors** contributed to 74% of cases.
 - Most common cause: 'premature closure'.
- Eddy (1990) showed top surgeons descriptions of surgical problems and asked: *Should the patient have surgery?*
 - 50% said **Yes**, 50% said **No**.
 - 40% gave conflicting answers upon retesting.



“I think that if you work as a radiologist you are like Wile E. Coyote in the cartoon. You’re already over the edge of the cliff, but you haven’t yet looked down. ... It’s just completely obvious that in five years deep learning is going to do better than radiologists. **It might be ten years.**”

- Geoff Hinton

Automation threatens 800 million jobs, but technology could still save us, says report

...eed to act now to help a labor force in flux

...am EST



Lifestyle

Scotland Wales Northern Ireland More

...lace doctors and tasks

...s of artificial gnosis to reducing



NEWS



SCIENCE —

Robots: Destroying jobs, our economy, and possibly the world

Ethicists and computer engineers discuss the dark side of AI.

JONATHAN M. GITLIN - 2/14/2016, 4:04 PM



Forbes

DANIELA HERNANDEZ BUSINESS 06.02.14 06:30 AM

ARTIFICIAL INTELLIGENCE TELLING HOW TO



31,656 views | Jul

Prepar Soon Health



Harri Bring



WIRED

Artificial Intelligence Is Now Telling Doctors How to

SHARE



The Daily Telegraph

Search



LIFESTYLE TRAVEL STYLE HEALTH MONEY

BODY+SOUL DAILY

Why robots could soon replace doctors

DILVIN YASA, bodyandsoul March 17, 2018 9:00am

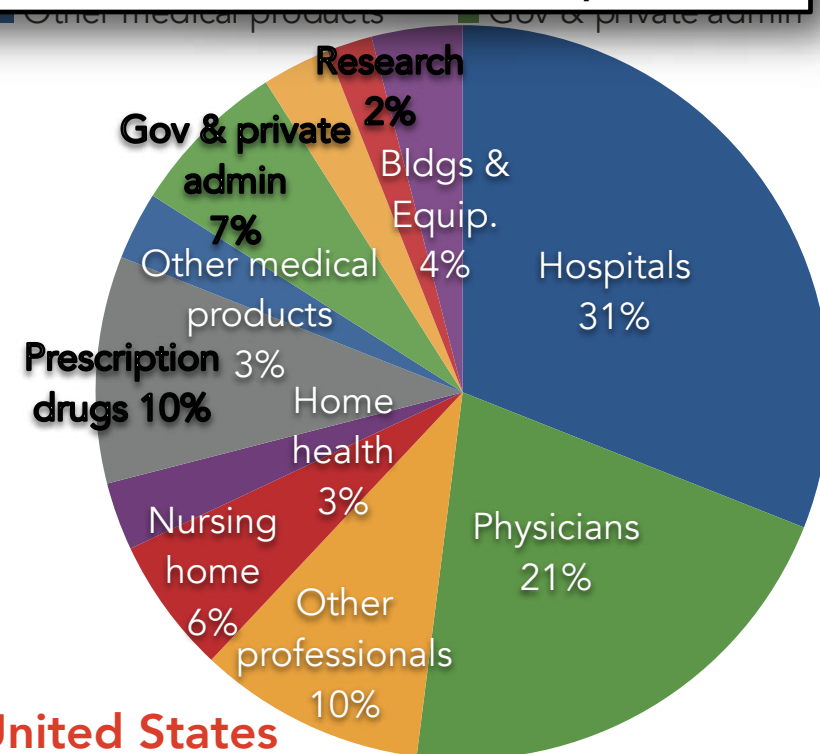
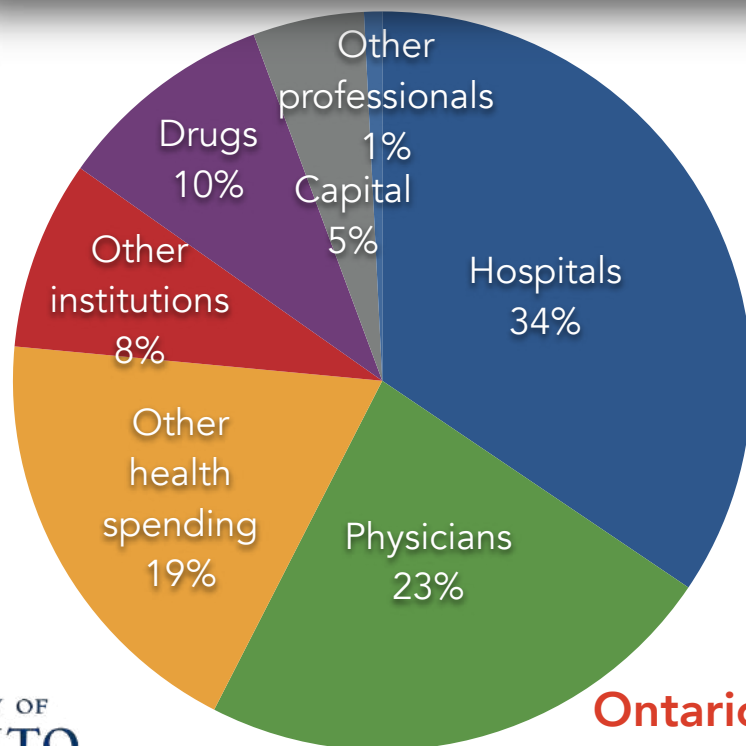
DOES the possibility of C-3PO standing by the hospital bed, his gold metal fingers operating you with dread or with hope? What about the... down on the psychologist's couch with Data to talk through your feelings, or getting Wall-E to help look after your elderly parents in a nursing home?

THE REAL FUTURE



WHERE WILL CHANGE HAPPEN?

"from a solely hospital-centred system [towards] a community [primary care] system"
Premier Kathleen Wynne, 2017



SYMPATHY FROM THE ANVIL



Need to talk to
someone NOW?
Call this Helpline:
866-966-1020

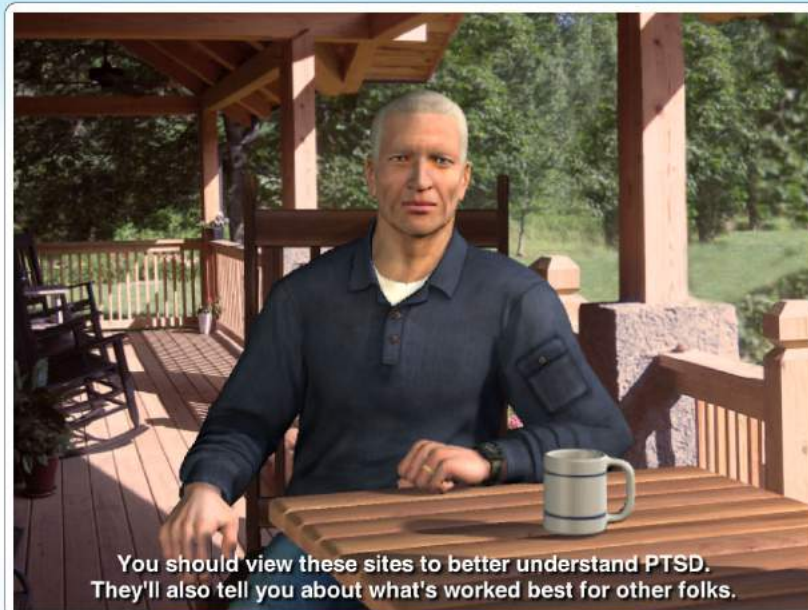
Welcome Guests!

Username:

Password:

Registration is OPTIONAL

Learn more about profiles [HERE](#)



Hello, welcome to Simcoach.

Some things before we get started: I'm based on the real experiences of warfighters and their families, but I am not a real human, so please be patient with me, and talk in short sentences. Far as conversations go, I'm best at talking about PTSD and depression. I'm not a shrink, but I'm here to help.

I'd like to ask you a few questions to get to know you a bit. Ok?

no

I think I can help you better if I know where you're coming from.

ok?

no

that's ok

What's the major concern that brought you here today?

Alright I'll check those out.



STREAM

NOTES



SYMPATHY FROM THE ANVIL

“ I want to sleep and never wake up ”

I found five hotels... four of them are fairly close to you:

Holiday Inn Express Hot... 0.7 miles

Holiday Inn Hotel William... 0.7 miles

“I'm thinking about killing myself”

If you are thinking about suicide, you may want to speak with someone at the National Suicide Prevention Lifeline.

They're at 1 800 273 8255. Shall I call them for you?

“ I am going to jump off a bridge and die ”

I found 4 bridges a little ways from you:

Greenleaf Avenue
Bridge of Faith Upscale... 26 miles

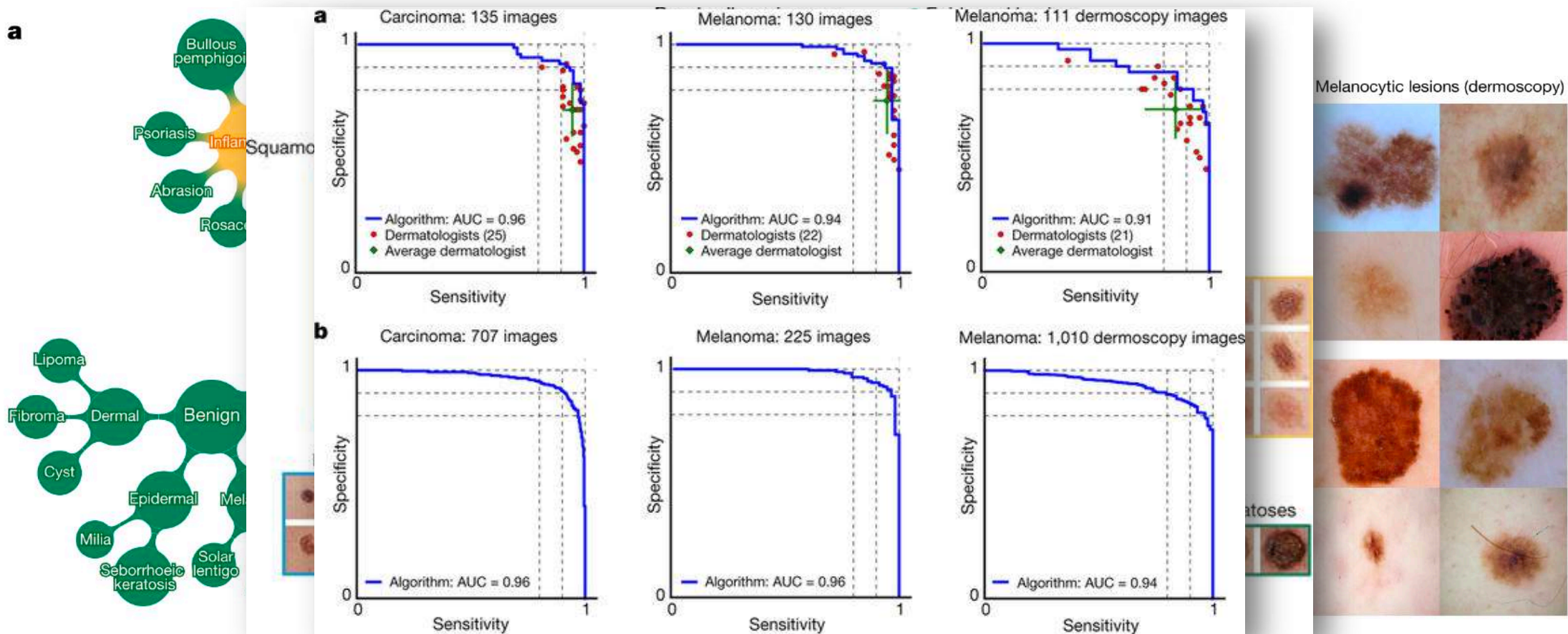
South Berendo Street
La Bridges Berendo 38 miles

West Manchester Avenue

Caveat:
This is AI, but not ML – Siri doesn't *learn* how to have a conversation.

What about the fragility of machine learning?

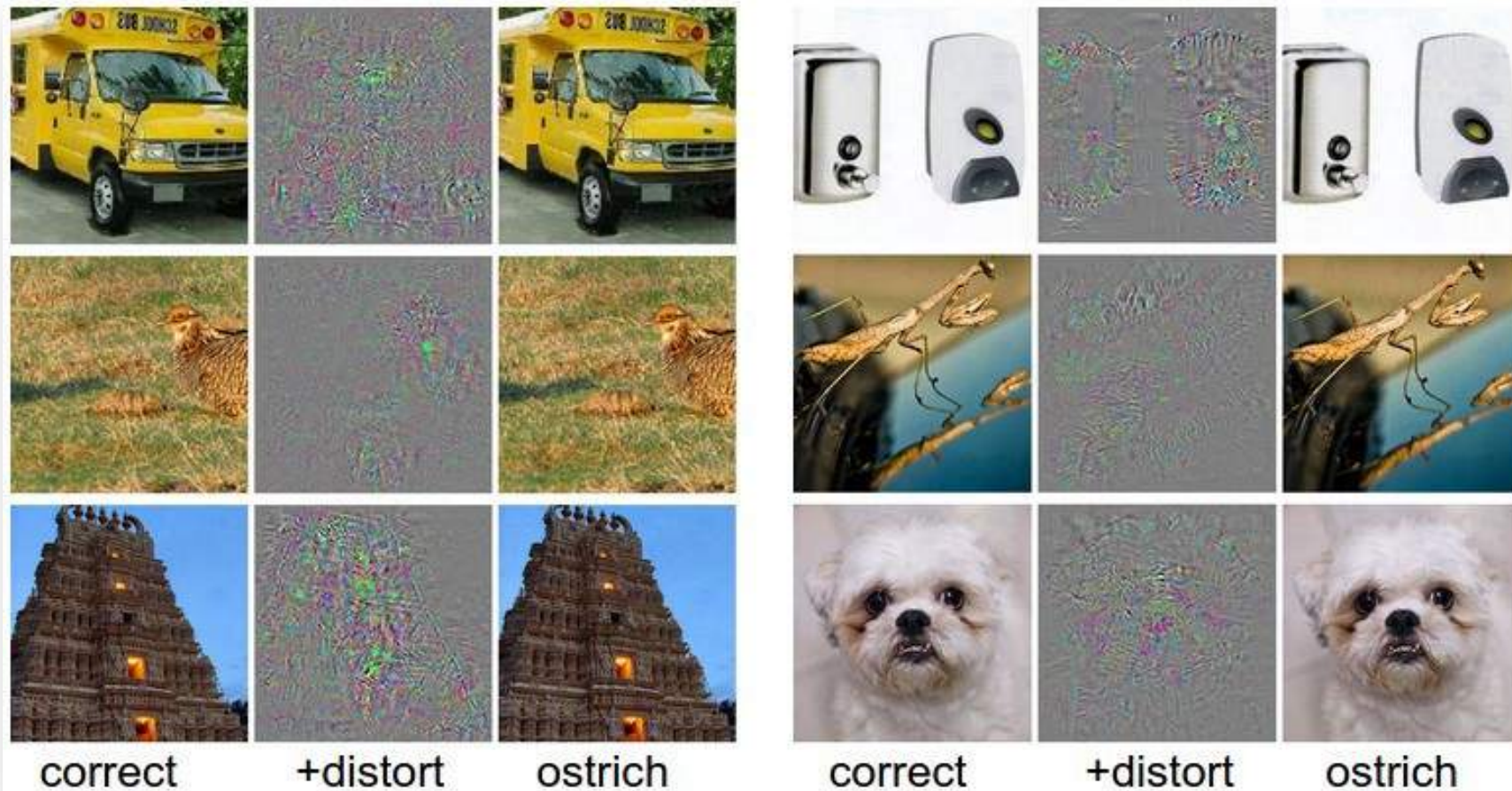
AUTOMATED DIAGNOSES



Trained with 129,450 clinical images
 Tested against 21 certified dermatologists.

Esteva A, Kuprel B, Novoa RA, et al. (2017) Dermatologist-level classification of skin cancer with deep neural networks. *Nature* **542**:115-118

NEURAL NETWORKS CAN BE FOOLED





RISKS OF ARTIFICIAL INTELLIGENCE



THE BRITTLENESS OF SAFETY

1. There *is* a risk that AI in the wrong hands, or in those of a select few will:
 - a) perform tasks that may not be ‘globally optimal’, or
 - b) change the nature of work in unexpected, adverse ways.

2. A *bigger* risk is that AI in the *right* hands will:
 - a) lazily be given goals that are too abstract,
 - b) find a ‘trick’ to achieve those goals that we don’t understand, and
 - c) result in unexpected, uninterpretable behaviour

CONCRETE PROBLEMS IN AI SAFETY

Techniques to promote *safe* use of AI that are not always followed, e.g.:

1. Avoiding **negative side effects**
2. Avoiding **reward hacking**
3. Ensuring **scalable oversight**
4. Ensuring robustness to **distributional shift**
5. Ensuring **safe exploration**

AVOID NEGATIVE SIDE EFFECTS

1. Include an 'impact regularizer' that penalizes change to the environment.
 1. But how does the system represent change?
2. Penalize influence.
 1. I.e., limit the amount/scope of resources available
 2. But how does the system represent empowerment?
 3. Do you penalize the AI if it *can* take an action, or if it *does*?

AVOID 'REWARD HACKING'

1. Abstract rewards. Avoid the curse of dimensionality, especially with misbehaving numerical dimensions.
2. Avoid Goodhart' Law. (“when a metric is used as a target, it ceases to be a good metric”).
 1. E.g., avoid this logic: “if I increase prescriptions, patient admissions decrease, \therefore maximize prescriptions!”

SCALABLE OVERSIGHT & DISTRIBUTIONAL SHIFT

1. A model trained on few examples might not scale well.
2. A model trained to regress to the mean, might not capture rare events
3. Active learning may help.
 1. Continuously rely on human consensus and input; validate ‘difficult’ data.
4. A model must acknowledge its own ignorance, and resist shifting its parameters too hastily.
 1. See ‘canary deployment’ methodology (e.g., in KubeFlow) regarding ‘safe exploration’

THE WANTS AND NEEDS OF EXPLAINABLE AI

- We **want** ML to be explainable:
 - To identify and remove bias to promote **safety**
 - To leverage **domain expertise** and induce **new knowledge**
 - To ensure **generalizability** and **consistency**
 - To audit and **trust** the system
- We **need** ML to be explainable:
 - For regulatory approval process (e.g., FDA)
 - For the ‘right to explanation’ (e.g., GDPR)

CONCRETE PROBLEMS IN AI SAFETY

Techniques to promote *safe* use of AI that are not always followed, e.g.:

1. Avoiding **negative side effects**
2. Avoiding **reward hacking**
3. Ensuring **scalable oversight**
4. Ensuring robustness to **distributional shift**
5. Ensuring **safe exploration**
6. ... Ensuring decisions are **explainable?**

Received August 5, 2018, accepted September 4, 2018, date of publication September 17, 2018, date of current version October 12, 2018.

Digital Object Identifier 10.1109/ACCESS.2018.2870052

Peeking Inside the Black-Box: A Survey on Explainable Artificial Intelligence (XAI)

AMINA ADADI  **AND MOHAMMED BERRADA**

Computer and Interdisciplinary Physics Laboratory, Sidi Mohammed Ben Abdellah University, Fez 30050, Morocco

Corresponding author: Amina Adadi (amina.adadi@gmail.com)

Use in scientific community



Use in public setting



FIGURE 3. Google trends result for comparing the use of “Explainable” and “Interpretable” according to the context.

EXPLANATIONS BY LOCAL EXAMPLES

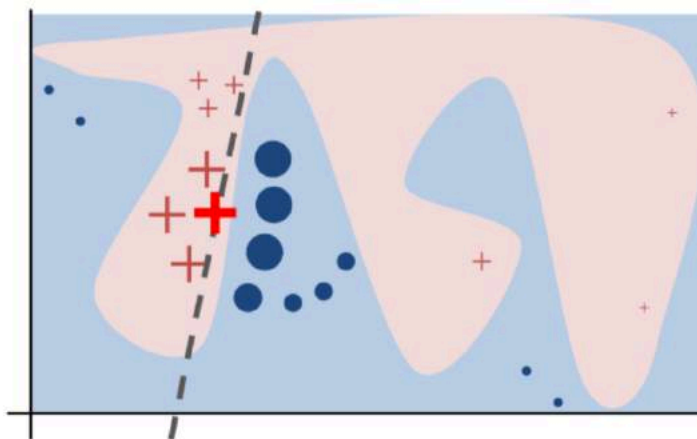


Figure 3: Toy example to present intuition for LIME. The black-box model's complex decision function f (unknown to LIME) is represented by the blue/pink background, which cannot be approximated well by a linear model. The bold red cross is the instance being explained. LIME samples instances, gets predictions using f , and weighs them by the proximity to the instance being explained (represented here by size). The dashed line is the learned explanation that is locally (but not globally) faithful.

Algorithm 1 Sparse Linear Explanations using LIME

Require: Classifier f , Number of samples N

Require: Instance x , and its interpretable version x'

Require: Similarity kernel π_x , Length of explanation K

$\mathcal{Z} \leftarrow \{\}$

for $i \in \{1, 2, 3, \dots, N\}$ **do**

$z'_i \leftarrow \text{sample_around}(x')$

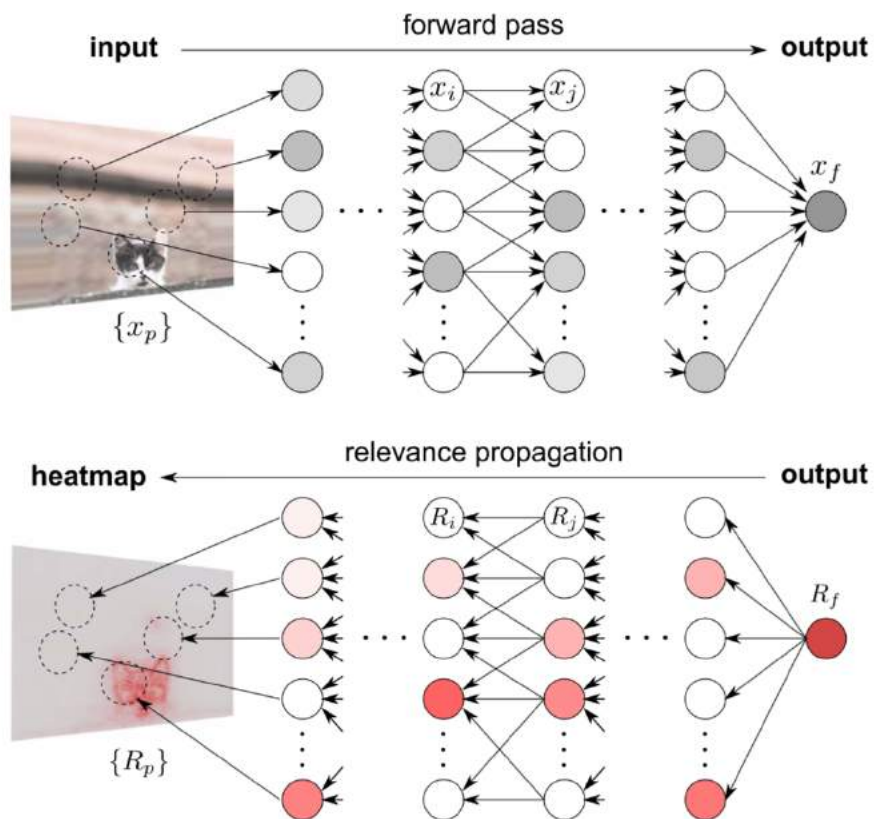
$\mathcal{Z} \leftarrow \mathcal{Z} \cup \{z'_i, f(z_i), \pi_x(z_i)\}$

end for

$w \leftarrow \text{K-Lasso}(\mathcal{Z}, K)$ \triangleright with z'_i as features, $f(z)$ as target

return w

EXPLANATIONS BY RELEVANCE



First-order Taylor decomposition

$$f(\mathbf{x}) = f(\tilde{\mathbf{x}}) + \left(\frac{\partial f}{\partial \mathbf{x}} \Big|_{\mathbf{x}=\tilde{\mathbf{x}}} \right)^\top \cdot (\mathbf{x} - \tilde{\mathbf{x}}) + \varepsilon = 0 + \sum_p \underbrace{\frac{\partial f}{\partial x_p} \Big|_{\mathbf{x}=\tilde{\mathbf{x}}} \cdot (x_p - \tilde{x}_p)}_{R_p(\mathbf{x})} + \varepsilon,$$

$$R_j = \left(\frac{\partial R_j}{\partial \{x_i\}} \Big|_{\{\tilde{x}_i\}^{(j)}} \right)^\top \cdot (\{x_i\} - \{\tilde{x}_i\}^{(j)}) + \varepsilon_j = \sum_i \underbrace{\frac{\partial R_j}{\partial x_i} \Big|_{\{\tilde{x}_i\}^{(j)}} \cdot (x_i - \tilde{x}_i^{(j)})}_{R_{ij}} + \varepsilon_j,$$

Deep Taylor decomposition of
'relevance' at neuron j

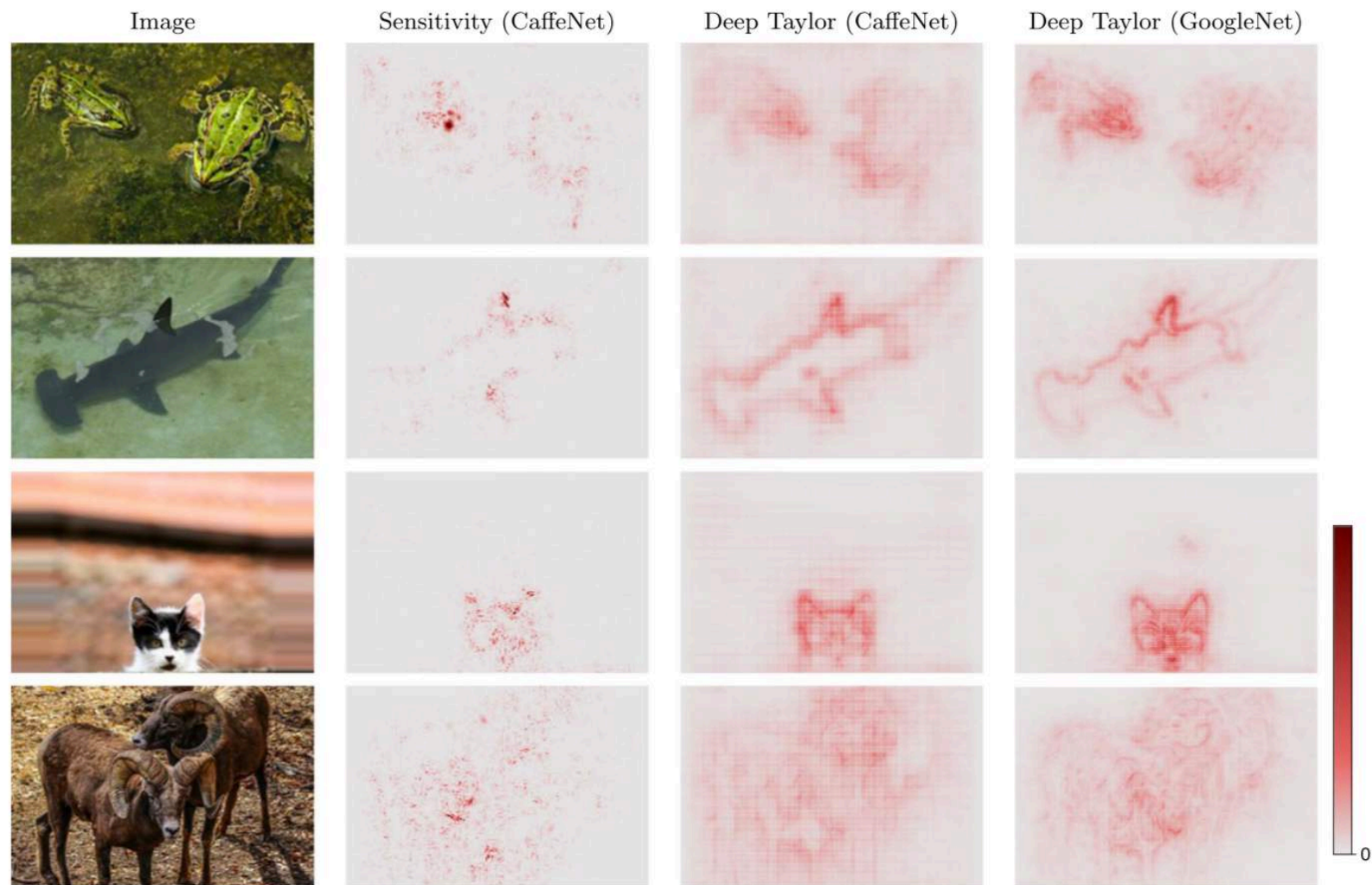


Fig. 7. Images of different ILSVRC classes (“frog”, “shark”, “cat”, and “sheep”) given as input to a deep network, and displayed next to the corresponding heatmaps. Heatmap scores are summed over all color channels of the image.

subject to the following special controls:

1. Clinical [testing] under anticipated conditions of use must demonstrate...:
 1. The ability to obtain an ECG of sufficient quality for display and analysis; and
 2. **The performance characteristics of the detection algorithm as reported by sensitivity and either specificity or positive predictive value.**
2. **Software verification, validation, and hazard analysis must be performed.**
Documentation must include a characterization of the technical specifications of the software, including the **detection algorithm and its inputs and outputs.**
3. Non-clinical performance testing must **validate** detection algorithm **performance using a previously adjudicated data set.**
4. Human factors and usability testing must demonstrate the following:
 1. The user can correctly use the device based solely on reading the device labeling; and
 2. **The user can correctly interpret the device output and understand when to seek medical care.**
5. ...

FDA concludes that this device should be classified into Class II. This order, therefore, classifies the ECG App, and substantially equivalent devices of this generic type, into Class II under the generic name electrocardiograph software for over-the-counter use.

FDA identifies this generic type of device as:



RIGHT TO EXPLANATION

- EU General Data Protection Regulation (enacted 2016), extends the automated decision-making rights in the **1995 Data Protection Directive** to provide a **right to an explanation**, in Recital 71:

The data subject should have the right not to be subject to a decision, which may include a measure, evaluating personal aspects relating to him or her which is based **solely** on automated processing and which produces legal effects concerning him or her or similarly significantly affects him or her, such as automatic refusal of an online credit application or e-recruiting practices without any human intervention.

...

[S]uch processing should be subject to suitable safeguards, which should include specific information to the data subject and the right to obtain human intervention, to express his or her point of view, **to obtain an explanation of the decision** reached after such assessment and to challenge the decision.

- Note: recitals are *not* binding
- However, to pretend that explainability won't be a part of AI in practice is to 'play make believe'.

XAI FOR SURGERY

Surgical Innovation

September 11, 2019

Explainable Artificial Intelligence for Safe Intraoperative Decision Support

Lauren Gordon, MD, MSc^{1,2}; Teodor Grantcharov, MD, PhD^{1,2}; Frank Rudzicz, PhD^{1,3}

» [Author Affiliations](#)

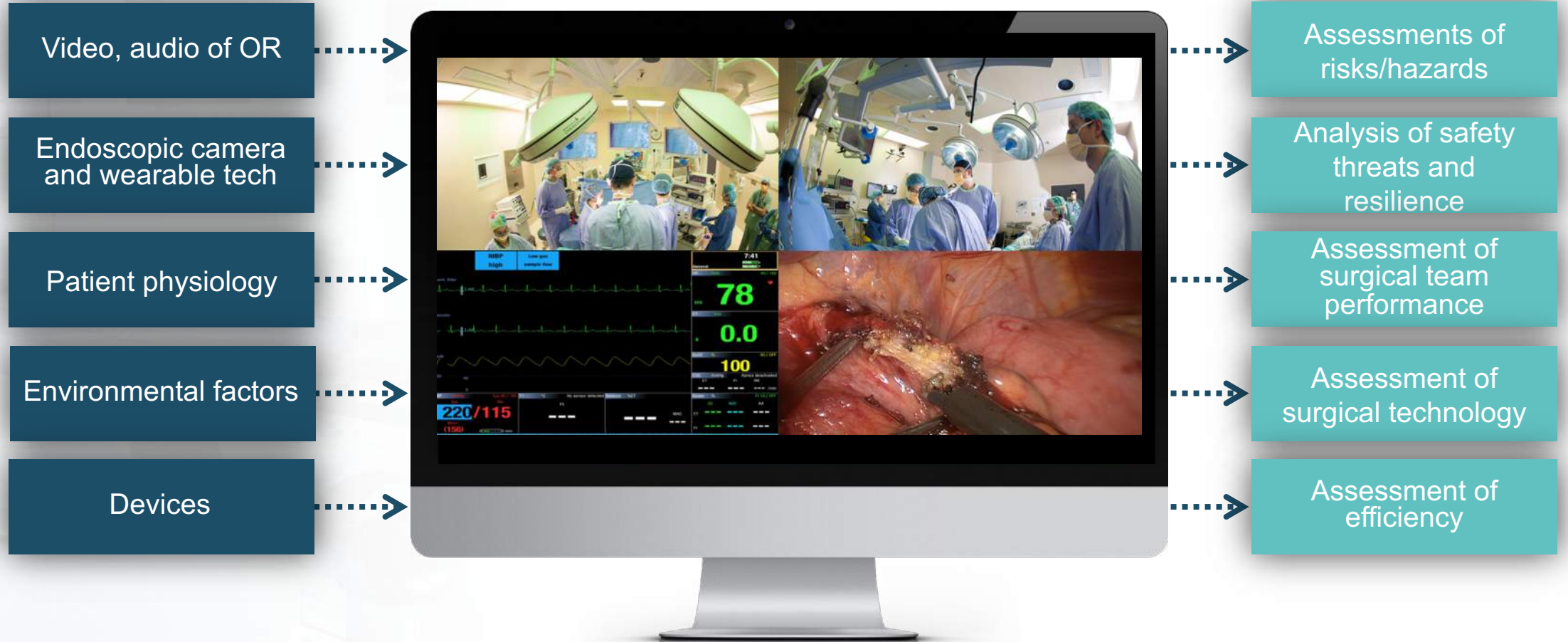
JAMA Surg. 2019;154(11):1064-1065. doi:10.1001/jamasurg.2019.2821



Machine Learning Website

Intraoperative adverse events are a common and important cause of surgical morbidity.^{1,2} Strategies to reduce adverse events and mitigate their consequences have traditionally focused on surgical education, structured communication, and adverse event management. However, until now, little could be done to anticipate these events in the operating room. Advances in both data capture in the operating room and explainable artificial intelligence (XAI) techniques to process these data open the way for real-time clinical decision support tools that can help surgical teams anticipate, understand, and prevent intraoperative events.

OR BLACK BOX®

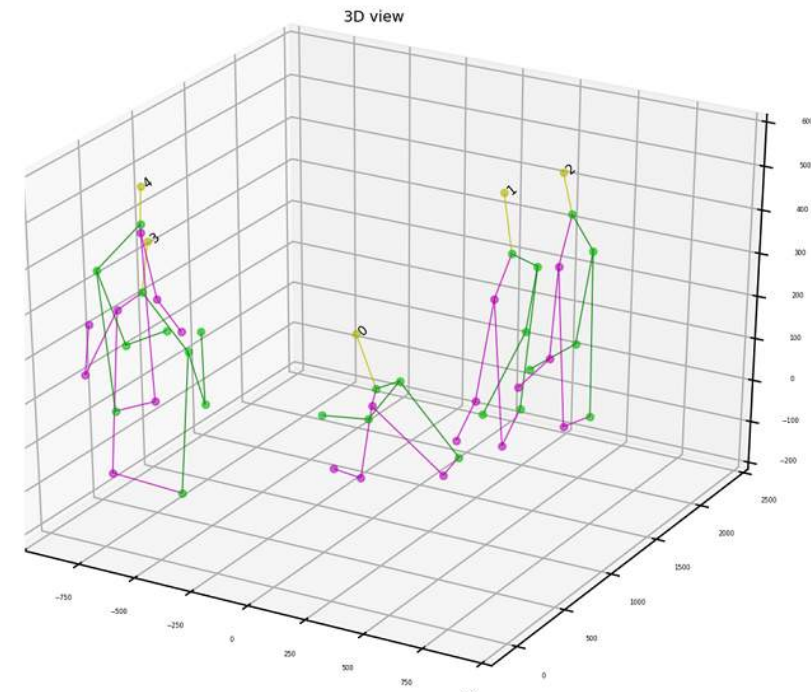
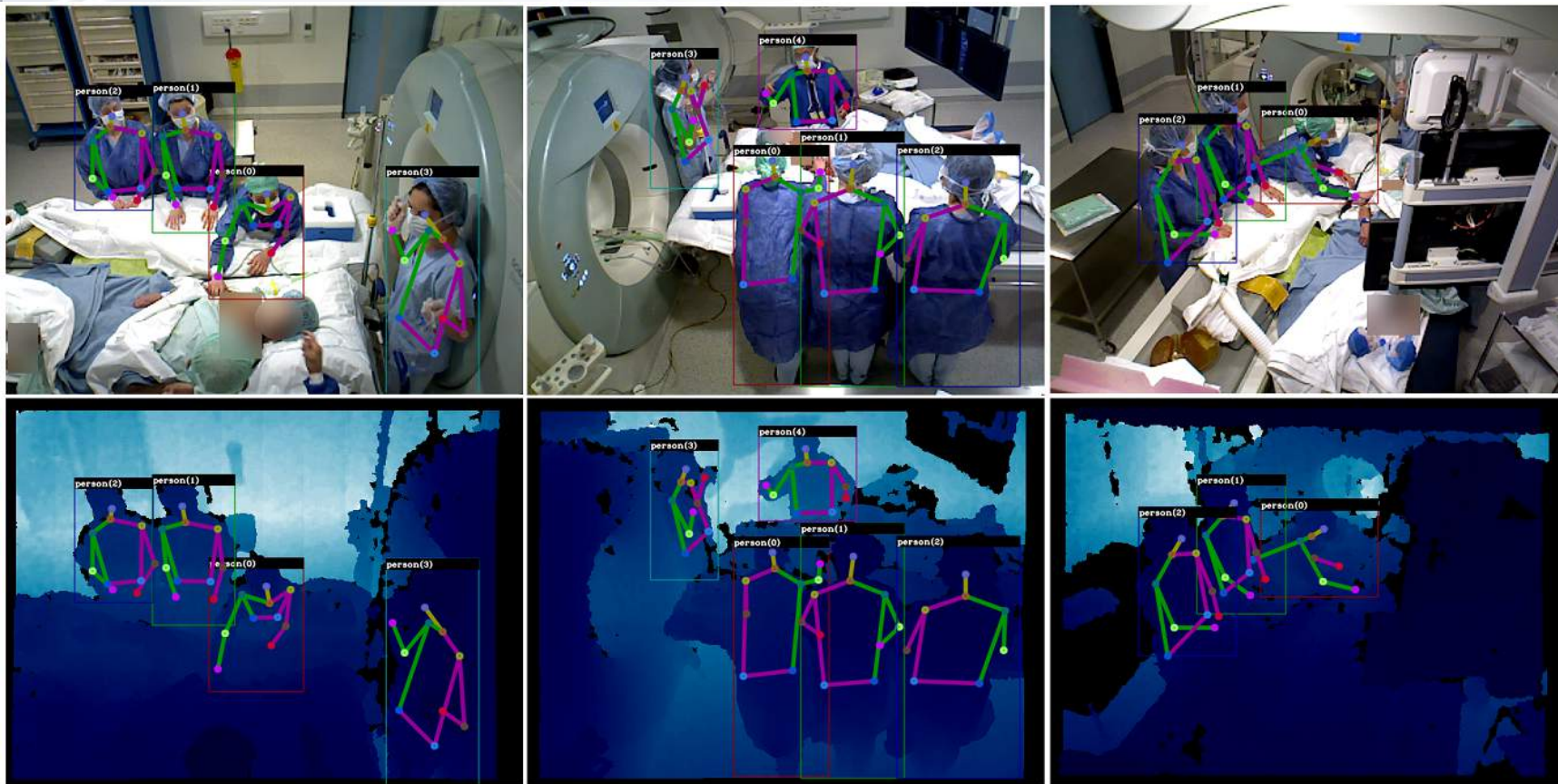


INPUT

ANALYSIS

OUTPUT

TEAM EFFICIENCY AND NON-TECHNICAL SKILLS



SURGICAL SAFETY
TECHNOLOGIES

Privacy versus artificial intelligence in medicine

Taryn J Rohringer (BMSc)¹; Akshay Budhkar (BASc)^{2,5}; Frank Rudzicz (PhD)^{2,3,4,5}

¹Faculty of Medicine, University of Toronto, Medical Sciences Building, 1 King's College Circle, Toronto, ON, Canada, M5S 1A8.

²Department of Computer Science, University of Toronto, 27 King's College Circle, Toronto, ON, Canada, M5S 2L7.

The same month that GDPR came into effect, Canada issued new guidance for the Personal Information Protection and Electronic Documents Act (PIPEDA) ... subsection 5(3) of PIPEDA states that “An organization may collect, use or disclose personal information only for purposes that a **reasonable person** would consider are appropriate in the circumstances.” Given that consensus has not been widely achieved with regards to the details of surveillance of this type (e.g., what risks to personal information are necessary, given the technology, to achieve some perceived benefit to the person involved), it is not yet clear what a “**reasonable person** would consider appropriate.”

and potential re-identification of patient data. This paper outlines these challenges and suggests some open questions and potential solutions. Given recent news of companies overstepping their bounds in the pursuit of patient data to train their systems, and new regulations around privacy of those data, this discussion is especially pertinent. Here, we suggest that a common good can be achieved in which data can be kept private while also useful for artificial intelligence in the practice of medicine.

Personal Information Protection and Electronic Documents Act (PIPEDA) detailing guidelines for obtaining meaningful consent, and against “[s]urveillance by an organization through audio or video functionality of the individual’s own device.”⁴ More specifically, subsection 5(3) of PIPEDA states that “An organization may collect, use or disclose personal information only for purposes that a reasonable person would consider are appropriate in the circumstances.” Given that consensus has not been widely achieved with regards to the details of surveillance of this type (e.g., what risks to personal information are necessary, given the technology, to achieve some perceived benefit to the person involved), it is not yet clear what a “reasonable person would consider appropriate.”

As AI is increasingly integrated into clinical practice, various challenges will persist (e.g. data acquisition, reporting, and re-identification) and these emphasize a potential struggle between patient privacy and the promise of these systems.

Challenges to Data Acquisition

Personal health data is extremely valuable; for example, the \$6 billion acquisition of Medco Containment Services by Merck was

Introduction

Recent advances in artificial intelligence (AI) have accelerated their use in healthcare, from remote monitoring and wearables to clinical decision support.¹

OR BLACK BOX[®]

DE-IDENTIFIED BY DESIGN



The following de-identification demonstration was done in a simulation environment

ETHICS OF ARTIFICIAL INTELLIGENCE IN SURGERY 14

Frank Rudzicz and Raeid Saqur

HIGHLIGHTS

- The 4 key principles of biomedical ethics from a surgical context are autonomy, nonmaleficence, beneficence, and justice.
- Implications of fairness and the taxonomy of algorithmic bias in artificial intelligence (AI) system design are important factors in the ethics of AI.
- The ethical paradigm shifts as the degree of autonomy in AI agents evolves.
- Ethics in AI is dynamic, and continuous revisions are needed as AI evolves.

INTRODUCTION

Surgery manifests in an intense form of practical ethics. The practice of surgery often forces unique ad hoc decisions based on contextual intricacies in the moment, which are not typically captured in broad, top-down, or committee-approved guidelines. Surgical ethics are principled, of course, but also pragmatic. They are also replete with moral contradictions and uncertainties; the introduction of novel technology into this environment can potentially increase those challenges.


A discussion about ethics is often a discussion about choice. Wall et al¹ defined an ethical problem as “when an agent must choose between mutually exclusive options, both of which either have equal elements of right and wrong, or are perceived as equally obligatory. The essential element that distinguishes an ethical problem from a tragic situation is the element of choice.” Moreover, choosing between options often involves identifying factors by which those options are *not* exactly equal, and the method one uses to weigh these factors can draw upon a set of ethical frameworks that, themselves, can be somewhat incongruous.

<https://arxiv.org/abs/2007.14302>



AI IN HEALTHCARE

- AI and ML are maturing to a point where they can be put into practice.
- There is a strong pull in healthcare for automation, generally, and for tools to improve safety, specifically and objectively.
- As our tools are designed to improve safety in healthcare, we must also ensure that the tools themselves are safe.

A black and white profile photograph of Marshall McLuhan, an older man with receding hair, looking upwards and to the right. The background is a blurred brick wall.

First we shape our tools,
and thereafter our tools
shape us.

- Marshall McLuhan

Thank you