

# When the digital Doctor should admit "I don't know"

June 24, 2020

## Abstract

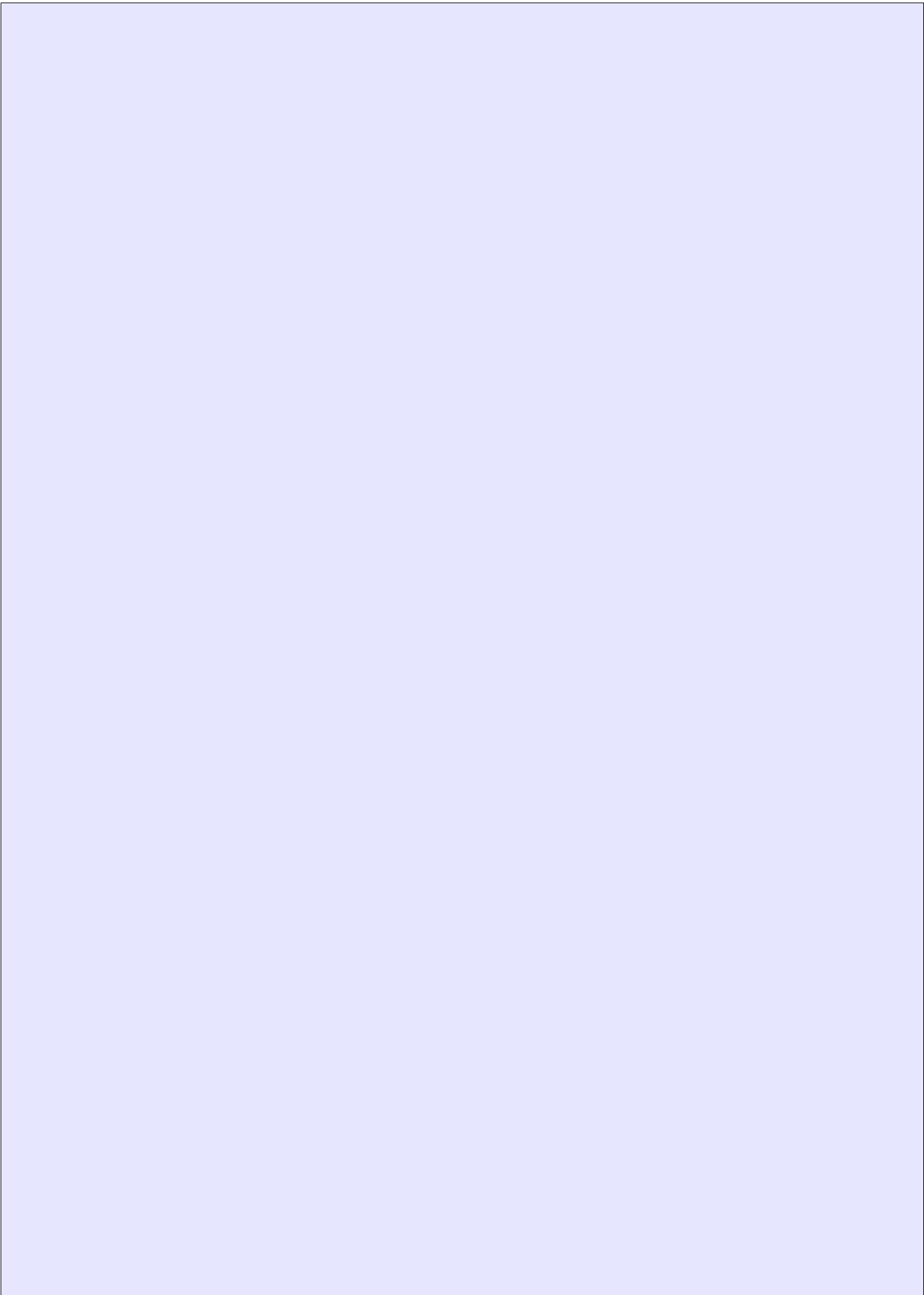
The meteoric rise of AI in general and Deep Learning in particular is generating great excitement throughout academia and commerce, and in particular in medicine[?, ?]. With some high-profile claims [] that AI will soon replace humans in many medical specialties.

In this position paper we present an alternative view. We contrast *Artificial Intelligence* with *Intelligence Augmentation* and argue that the second is more likely to benefit the patient than the first. We provide evidence to this argument and present a vision in which easier decisions are delegated to computers, while the more difficult ones are handled by humans.

## Introduction

Digital technology is causing a sea-change in all parts of the medical profession. In particular the meteoric rise of AI in general and deep learning in particular raises the possibility that doctors will be replaced computers [Mukherjee(2017)]. The father of deep learning, Geoff Hinton, said in 2017: "It's just completely obvious that that in ten years deep learning is going to do better than Radiologists ... They should stop training radiologists now".

Other deep learning researchers provide a more nuanced perspective. Sebastian Thrun [Mukherjee(2017), Esteva et al.(2017)Esteva, Kuprel, Novoa, Ko, Swetter, Blau, and Thrun] argues that "... deep learning devices will not replace dermatologists and radiologists. They will *augment* professionals, offering the expertise and assistance".



**What would IA look like when applied to medicine?** that is the question we aim to answer here. We argue that an important ingredient of the answer is to introduce to AI agents a level of humility. Specifically, to design classifiers, such as DNNs, to say "I don't know".

## **Labels, ground truth and testing**

Roughly speaking, machine learning (ML) can be divided into *un-supervised* learning and *supervised* learning. In both, the task of the learning algorithm is transform a set of *examples* into a *model*. In un-supervised learning the



Skin cancer diagnosis using Deep Neural NetworksSkin cancer diagnosis using Deep Neural

One  
of  
the  
pa-  
pers  
that  
pro-  
vided

ev-  
i-  
dence  
that  
deep  
neu-  
ral  
net-  
works

might  
be  
able

to  
out-  
per-  
form

hu-  
mans

is  
the  
work

of  
Es-  
teva

et

al [Esteva et al.(2017)Esteva, Kuprel, Novoa, Ko, Swetter, Blau, and Thrun].

They  
trained

a  
Deep

neu-  
ral  
net-  
work

to  
clas-  
sify

im-  
ages  
of

skin  
into

three  
cat-

e-  
gories:

be-

In their famous work, Esteva et al. set out to show that a classifier trained by machine learning can perform as well as or better than expert dermatologists. In this application of supervised learning each example consists of an input image of a skin patch and an output label that is “benign” or “malignant”

As they wanted to compare the system to human dermatologists they needed a better ground truth than that provided by the dermatologists. To that end they used the diagnosis of a biopsy as ground truth. It is arguable that this label is more accurate than the one given by the dermatologist, even though it depends on the human judgement of the pathologist.

However, even if we assume that pathologists labels are more reliable than dermatology labels, the requirement that each example corresponds to a biopsy introduces a significant bias. Under normal circumstances, patients get biopsied only if the dermatologist thinks there is a chance of **malignancy**. Therefore, the set of biopsied examples is biased towards **malignancy**. It is likely that using a classifier trained in this way on an unfiltered stream of patients will increase the number of patients unnecessarily getting a biopsy.

## Uncertainty in medicine

Medicine is rife with risk and uncertainty. An incorrect diagnosis or treatment can cost the patient his life and the doctor her license.

Uncertainty has many causes, we discuss some of those below.

Medical  
de-  
vices  
use  
a  
va-  
ri-  
ety  
of  
bio-  
sensors  
that  
record,  
dis-  
play  
and  
dis-  
tribute  
dif-  
fer-  
ent  
bio-  
met-  
rics,  
rang-  
ing  
from  
vi-  
tal  
signs  
such  
as  
heart  
rate,  
oxy-  
gen  
sat-  
u-  
ra-  
tion  
and  
blood  
pres-  
sure,  
to  
high-  
frequency  
wave-  
forms  
such  
as  
ECG,  
EEG,  
res-  
pi-  
ra-





The  
Amer-  
i-  
can  
Academy  
of  
Sleep  
Medicine  
(AASM)  
pub-  
lishes  
cri-  
te-  
ria  
for  
man-  
ual  
sleep  
stage  
and  
sleep  
ap-  
nea  
an-  
no-  
ta-  
tion  
from  
the  
gold  
stan-  
dard  
sleep  
study  
in-  
stru-  
ment,  
the  
polysomno-  
gram  
(PSG).  
This  
an-  
no-  
ta-  
tion  
is  
based  
on  
man-  
ual  
anal-  
y-  
sis  
of  
biosig-



+Cohen's  
kappa,  
Yoav  
:  
In-  
tro-  
duce  
Hoen  
kappa.  
In-  
stead  
of  
math-  
e-  
mat-  
i-  
cal  
def-  
i-  
ni-  
tion,  
I  
would  
in-  
ter-  
pret  
some  
spe-  
cific  
val-  
ues:  
the  
value  
that  
cor-  
responds

Quantification of inter-rater agreement rateQuantification of inter-rater agreement rate

to  
per-  
fect  
agree-  
ment  
for  
pos-  
i-  
tive  
per-  
fect  
agree-  
ment  
for  
neg-  
a-  
tive,  
the  
val-  
ues



Urodynamic studies provide the best bladder and sphincter functional data for urologists to decide how to treat patients at risk for renal damage [Abrams(2003)]. While it has been extensively studied and applied in clinics, the main issue that plagues

Medical uncertainty as manifest low inter-rater agreement consequence, can be found in many clinical problems (see blocks on ...)For example, the low agreement might come from the “extrapolation error”; that is, when we apply the developed protocol to the population different from the population that we collect the evidence for the protocol [Brosnan et al.(2015)Brosnan, La Gerche, Kumar, Lo, Kalman, and Prior]. In other situation, the variability among subjects is so big that it limits the development of a more quantitative protocol [Venhola et al.(2003)Venhola, Reunanen, Taskinen, Lahdes-Vasama, and Uhari]. In some situations, when the needed information is missing, it is challenging to make a differential diagnosis [Moncada et al.(2011)Moncada, Rueda, Maci

A direct consequence of the low inter-rater agreement rate is that the trained intelligent system might be questionable. It is clear that such intelligent system is questionable and might raise concerns. Recently, various regulations in this regard have been proposed [Price and Nicholson(2014), Ford et al.(2016)Ford, Price, and Nicholson].

**Hautieng : should we jump into GDPR? Yoav : what is GDPR?**

Now, suppose we are able to eliminate all challenges from data calibration and validation issues, and we can provide as much information as possible to train the intelligence system. Even under this assumption, it is clear that the system still suffers from the protocol limitation or knowledge gap issues. Can such system be useful in clinics? To answer this question, we should not forget that physicians also follow the same protocol and have knowledge gaps. Depending on the clinical problems, and the experience of physicians under consideration, the agreement rate varies. Usually, intern doctors know the least, while a senior attending knows the most. It is natural that we trust a senior expert more, but it does not mean that we do not trust a junior intern doctor.

In general, diagnosis is performed by comparing observed symptoms to past experience. Roughly speaking there are two ways to make this comparison: *recognition* and *elimination*. Recognition is a fast, typically non-





A  
di-  
rect  
con-  
se-  
quence  
of  
this  
low  
inter-  
rater  
agree-  
ment  
is  
a  
ques-  
tion-  
able  
trained  
“ar-  
ti-  
fi-  
cial  
in-  
tel-  
li-  
gence”.  
It  
is  
pos-  
si-  
ble  
that  
we  
mag-  
i-  
cally  
ob-  
tain  
a  
dataset  
that  
con-  
tains  
in-  
for-  
ma-  
tion  
that  
is  
suf-  
fi-  
cient  
for  
the  
de-

## Sources of uncertainty in medical diagnosis.

- **The diagnostic process of elimination**
- **Data Quality, Calibration, resolution** Discuss issue as placement of sensors, .

## Hiding Uncertainty

- **Psychological reasons** Both doctor and patient prefer the projection of certitude.
- **Protocols** –done
- **diagnostic devices** Secrecy of the internal code limits the trustworthiness of the alarms.–done
- **Alarm Fatigue**–done

How to quantify IDK? We should discuss how to quantify the confidence, or certainty, a physician has when making a decision. Clearly, experience leads to confidence. With more experience aggregated, diagnostic options that contradict the accumulated experience are eliminated, and hence more problems that need to be handled by the elimination process can be handled by the recognition process. However, facing our complicated human body, it is almost not possible for any single physician to aggregate all necessary experience to be confident about anything, so IDK is still an option. A practical and simple way to increase diagnostic certainty is to solicit the experience of a diverse group of doctors via discussion. If there is a clear majority for one diagnostic outcome, then the overall confidence in that diagnostics is high. While this voting procedure might be guarantee the optimal outcome, it eliminates the uncertainty during the whole procedure. With this certain procedure, even if the outcome is negative, it can be traced back and accumulate evidence and experience.

This certainty is very different from the the conditional probability of the disease given the diagnostic. The first is akin to saying: 95% of the dermatologists would give the same

Certainty and conditional probabilityCertainty and conditional probability

diagnostics. The second defines

## Uncertainty in Machine Learning

One can define “confidence” in machine learning. The definition follows a similar logic to the one used for human diagnosticians in the previous section. The yardstick by which we measure confidence of predicting a label is “how much do alternative labels contradict previous experience?”. More formally, we ask how much do we need to change the training data so that it supports an alternative label.

Uncertainty vs. accuracy using ROC curves

- Bootstrap samples.
- Samples from different hospitals.
- Easy and hard cases.

## Human decisions and Intelligence augmentation

Computers are an integral part of medical practice. From electronical medical records to medical instrumentation to billing, hospitals and clinics cannot function without computers. By some measures computers can already make better diagnosis than human doctors. The question is not *whether* computer diagnostics will become part of medical practice, the question is *how*.

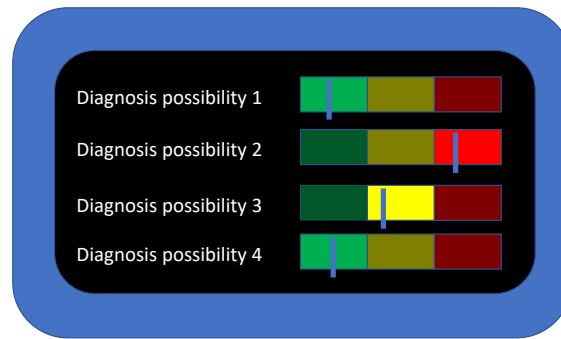
Some claim that human doctors and nurses are heading to extinction, following the fate of manufacturing jobs and bank cashiers. Our prediction is that computers will change the nature of medical work, but that it will increase, rather than decrease, the number of healthcare workers, especially in the care of chronic disease and aging.

We believe computers *can* perform accurate diagnosis for cases where different doctors are likely to agree. In other cases which are diagnostic gray area the computer will output “I don’t know” and transfer the responsibility to the doctor. In most cases, the doctor cannot say “I don’t know” because she is responsible for the patients health. On the other hand, resolving the diagnostic question is not her only choice. She can consult another doctor or the literature, ask for additional tests, or decide on a treatment based on available information. Deciding between these options requires much more than diagnostic information. It involves understanding the patient’s emotional, mental and financial state, the patient’s support system, the strengths and weaknesses of the hospital in which this is taking place etc.

Over time, computers will be able to take into consideration more and more of this complex information. However, for the foreseeable future, it is unlikely that computers will be given the responsibility to make medical *decisions*. Computers will take on much of the diagnostics and alarm tasks, improving the accuracy and timeliness of the doctors actions. Computers will output IDK in gray areas and will leave the decision making to the human doctor. Giving the computer the authority to make decisions currently done by human doctors will deprive the patient the human attention of the doctor.

Some of the digitization of the medicine has come between patients and doctors. The need to record all activity into EMR system require doctors to spend more time at the keyboard, reducing the amount of time of physical examination and discussion. We believe that IA can move medicine in the opposite direction, letting the computer make the common noncontroversial diagnostics and giving the patient more time to interact with the patient.

For IA technology to be widely adopted, the nurses and doctors that use them should experience an improvement in their practice. Suppose that the display of the diagnostics computer uses a three color code for each . Green indicates a confident negative diagnostic, red corresponds to a confident positive diagnosis. Finally, yellow corresponds to IDK, meaning that the computer cannot confirm or reject the diagnostic outcome.



The thresholds which define the three ranges ....

We finish this section with a few application areas which seem ready for applications of IA.

- **Computer aided diagnostics for large-scale data**

Medical imaging devices such as digital X-ray, CT, EMR and scanning microscope generate many gigabytes of data for each patient. Radiologists and pathologists spend their days analyzing these images to diagnose the patient. The large size and high resolution of the images on the one hand, and the time limitation on the analyst on the other imply that the analyst has to quickly narrow down the suspicious region, increase the chance of missing dangerous abnormalities.

IA can help the pathologist by suggesting locations in the high resolution image that might contain cancer nodules [1].

directing her attention to the parts of the image that are

- **Adaptive Patient monitors**

- **Dissemination of expertise** Computers, trained by experts, can help novices. Serves a function similar to score-cards.

Teaching young diagnostics

## Summary

## References

- [Mukherjee(2017)] Siddhartha Mukherjee. A.i. versus m.d.: What happens when diagnosis is automated? *The New Yorker*, April 2017.
- [Esteva et al.(2017)] Esteva, Kuprel, Novoa, Ko, Swetter, Blau, and Thrun] Andre Esteva, Brett Kuprel, Roberto A Novoa, Justin Ko, Susan M Swetter, Helen M Blau, and Sebastian Thrun. Dermatologist-level classification of skin cancer with deep neural networks. *Nature*, 542(7639):115–118, 2017.
- [Ashby(1957)] W Ross Ashby. An introduction to cybernetics. 1957.
- [Engelbart(1962)] Douglas C Engelbart. Augmenting human intellect: A conceptual framework. *Menlo Park, CA*, 1962.
- [Mehta et al.(2011)] Mehta, Granton, Lapinsky, Newton, Bandayrel, Little, Siau, Cook, Ayers, Singer, et al.] Sangeeta Mehta, John Granton, Stephen E Lapinsky, Gary Newton, Kristofer Bandayrel, Anjuli Little, Chuin

- Siau, Deborah J Cook, Dieter Ayers, Joel Singer, et al. Agreement in electrocardiogram interpretation in patients with septic shock. *Critical care medicine*, 39(9):2080–2086, 2011.
- [Atiya et al.(2003)Atiya, Kurth, Berger, Buring, and Kase] Monika Atiya, Tobias Kurth, Klaus Berger, Julie E Buring, and Carlos S Kase. Interobserver agreement in the classification of stroke in the women’s health study. *Stroke*, 34(2):565–567, 2003.
- [Johnson et al.(2016)Johnson, Ghassemi, Nemati, Niehaus, Clifton, and Clifford] Alistair E W Johnson, Mohamad M. Ghassemi, Shamim Nemati, Katherine E. Niehaus, David Clifton, and Gari D. Clifford. Machine Learning and Decision Support in Critical Care. *Proceedings of the IEEE*, 104(2):444–466, 2016.
- [Fleischman et al.(2019)Fleischman, Ciliberto, Rozanski, Parwani, and Bernstein] William Fleischman, Bethany Ciliberto, Nicole Rozanski, Vivek Parwani, and Steven L Bernstein. Emergency department monitor alarms rarely change clinical management: An observational study. *The American journal of emergency medicine*, 2019.
- [Cvach(2012)] Maria Cvach. Monitor alarm fatigue: an integrative review. *Biomedical instrumentation & technology*, 46(4):268–277, 2012.
- [BRIEF(2020)] EXECUTIVE BRIEF. Top 10 health technology hazards for 2020. 2020. URL <https://www.ecri.org/landing-top-10-patient-safety-concerns-2020>. Top 10 patient safety concerns of 2020 from ECRI Institute.
- [Behar et al.(2013)Behar, Oster, Li, and Clifford] Joachim Behar, Julien Oster, Qiao Li, and Gari D Clifford. Ecg signal quality during arrhythmia and its application to false alarm reduction. *IEEE transactions on biomedical engineering*, 60(6):1660–1666, 2013. alarm fatigue prevention for arrhythmia.
- [Bai et al.(2016)Bai, Do, Ding, Palacios, Shahriari, Pelter, Boyle, Fidler, and Hu] Yong Bai, Duc Do, Quan Ding, Jorge Arroyo Palacios, Yalda Shahriari, Michele M Pelter, Noel Boyle, Richard Fidler, and Xiao Hu. Is the sequence of superalarm triggers more predictive than sequence of the currently utilized patient monitor alarms? *IEEE Transactions on Biomedical Engineering*, 64(5):1023–1032, 2016. a new alarm system called superalarm to avoid alarm fatigue.
- [Zong et al.(2016)Zong, Nielsen, Gross, Brea, and Frassica] Wei Zong, Larry Nielsen, Brian Gross, Juan Brea, and Joseph Frassica. A practical algorithm to reduce false critical ecg alarms using arterial blood pressure and/or photoplethysmogram waveforms. *Physiological measurement*, 37(8):1355, 2016.
- [Winters et al.(2018)Winters, Cvach, Bonafide, Hu, Konkani, O’Connor, Rothschild, Selby, Pelter, McLean, et al.] Bradford D Winters, Maria M Cvach, Christopher P Bonafide, Xiao Hu, Avinash Konkani, Michael F O’Connor, Jeffrey M Rothschild, Nicholas M Selby, Michele M Pelter, Barbara McLean, et al. Technological distractions (part 2): a summary of approaches to manage clinical alarms with intent to reduce alarm fatigue. *Critical care medicine*, 46(1):130–137, 2018.
- [Hu(2019)] Xiao Hu. An algorithm strategy for precise patient monitoring in a connected healthcare enterprise. *NPJ digital medicine*, 2(1):1–5, 2019.
- [Komorowski et al.(2018)Komorowski, Celi, Badawi, Gordon, and Faisal] Matthieu Komorowski, Leo A Celi, Omar Badawi, Anthony C Gordon, and A Aldo Faisal. The artificial intelligence clinician learns optimal treatment strategies for sepsis in intensive care. *Nature medicine*, 24(11):1716–1720, 2018. sepsis treatment via AI.

- [Huddar et al.(2016)Huddar, Desiraju, Rajan, Bhattacharya, Roy, and Reddy] Vijay Huddar, Babu Koundinya Desiraju, Vaibhav Rajan, Sakyajit Bhattacharya, Shourya Roy, and Chandan K Reddy. Predicting complications in critical care using heterogeneous clinical data. *IEEE Access*, 4:7988–8001, 2016. ML for complication prediction in ICU.
- [Meyer et al.(2018)Meyer, Zverinski, Pfahringer, Kempfert, Kuehne, Sündermann, Stamm, Hofmann, Falk, and Eickhoff] Alexander Meyer, Dina Zverinski, Boris Pfahringer, Jörg Kempfert, Titus Kuehne, Simon H Sündermann, Christof Stamm, Thomas Hofmann, Volkmar Falk, and Carsten Eickhoff. Machine learning for real-time prediction of complications in critical care: a retrospective study. *The Lancet Respiratory Medicine*, 6(12): 905–914, 2018. ML for complication prediction in ICU.
- [Ruskin and Shelley(2005)] Keith J. Ruskin and Kirk H. Shelley. Patent medicine and the "black box". *Anesthesia and Analgesia*, 100(5):1361–1362, 2005. ISSN 00032999.
- [Feldman(2006)] Jeffrey M. Feldman. Can clinical monitors be used as scientific instruments? *Anesthesia and Analgesia*, 103(5):1071–1072, 2006.
- [Shelley and Barker(2016)] Kirk H. Shelley and Steven J. Barker. Disclosures, what is necessary and sufficient? *Anesthesia and Analgesia*, 122(2):307–308, 2016.
- [Cannesson and Shafer(2016)] Maxime Cannesson and Steven L. Shafer. All boxes are black. *Anesthesia and Analgesia*, 122(2):309–317, 2016.
- [Lin et al.(2019)Lin, Lo, Lin, Frasch, and Wu] Yu-Ting Lin, Yu-Lun Lo, Chen-Yun Lin, Martin G Frasch, and Hau-Tieng Wu. Unexpected sawtooth artifact in beat-to-beat pulse transit time measured from patient monitor data. *PloS one*, 14(9), 2019.
- [Gesche et al.(2012)Gesche, Grosskurth, Kuchler, and Patzak] Heiko Gesche, Detlef Grosskurth, Gert Kuchler, and Andreas Patzak. Continuous blood pressure measurement by using the pulse transit time: comparison to a cuff-based method. *European journal of applied physiology*, 112(1):309–315, 2012.
- [Iber et al.(2007)Iber, Ancoli-Isreal, Jr., , and Quan] C. Iber, S. Ancoli-Isreal, A. Chesson Jr., , and S. Quan. *The AASM Manual for Scoring of Sleep and Associated Events-Rules: Terminology and Technical Specification*. American Academy of Sleep Medicine, 2007.
- [Berry et al.(2012)Berry, Budhiraja, and et al.] R. B. Berry, D. G. Budhiraja, and et al. Rules for scoring respiratory events in sleep: update of the 2007 AASM manual for the scoring of sleep and associated events. *J Clin Sleep Med*, 8(5):597–619, 2012.
- [Norman et al.(2000)Norman, Pal, Stewart, Walsleben, and Rapoport] Robert G Norman, Ivan Pal, Chip Stewart, Joyce A Walsleben, and David M Rapoport. Interobserver agreement among sleep scorers from different centers in a large dataset. *Sleep*, 23(7):901–908, 2000.
- [Abrams(2003)] Paul Abrams. Describing bladder storage function: overactive bladder syndrome and detrusor overactivity. *Urology*, 62(5):28–37, 2003.
- [Venhola et al.(2003)Venhola, Reunanen, Taskinen, Lahdes-Vasama, and Uhari] Mika Venhola, Mikko Reunanen, Seppo Taskinen, Tuija Lahdes-Vasama, and Matti Uhari. Interobserver and intra-observer agreement in interpreting urodynamic measurements in children. *The Journal of urology*, 169(6):2344–2346, 2003.
- [Dudley et al.(2018)Dudley, Adams, Brock III, Clayton, Joseph, Koh, Merguerian, Pope IV, Routh, Thomas, et al.] Anne G Dudley, Mark C Adams, John W Brock III, Douglass B Clayton, David B Joseph, Chester J Koh,



- Paul A Merguerian, John C Pope IV, Jonathan C Routh, John C Thomas, et al. Interrater reliability in interpretation of neuropathic pediatric urodynamic tracings: an expanded multicenter study. *The Journal of urology*, 199(5):1337–1343, 2018.
- [Brosnan et al.(2015)Brosnan, La Gerche, Kumar, Lo, Kalman, and Prior] Maria Brosnan, Andre La Gerche, Saurabh Kumar, Wilson Lo, Jonathan Kalman, and David Prior. Modest agreement in ecg interpretation limits the application of ecg screening in young athletes. *Heart Rhythm*, 12(1):130–136, 2015.
- [Moncada et al.(2011)Moncada, Rueda, Macías, Suárez, Ortega, and Vélez] Diana Carolina Moncada, Zulma Vanessa Rueda, Antonio Macías, Tatiana Suárez, Héctor Ortega, and Lázaro Agustín Vélez. Reading and interpretation of chest x-ray in adults with community-acquired pneumonia. *The Brazilian Journal of Infectious Diseases*, 15(6):540–546, 2011.
- [Price and Nicholson(2014)] W Price and II Nicholson. Black-box medicine. *Harv. JL & Tech.*, 28:419, 2014.
- [Ford et al.(2016)Ford, Price, and Nicholson] Roger Allan Ford, W Price, and II Nicholson. Privacy and accountability in black-box medicine. *Mich. Telecomm. & Tech. L. Rev.*, 23:1, 2016.
- [Vordermark II(2019)] Jonathan S Vordermark II. *An Introduction to Medical Decision-Making: Practical Insights and Approaches*. Springer Nature, 2019.
- [Michel(2020)] Jeffrey B Michel. Thinking fast and slow in medicine. In *Baylor University Medical Center Proceedings*, volume 33, pages 123–125. Taylor & Francis, 2020.