

1 **The promises and pitfalls of reinforcement learning in healthcare**

2 Omer Gottesman^{*1}, Fredrik Johansson^{*2}, Matthieu Komorowski^{3,4}, Aldo Faisal⁵, David Sontag²,

3 Finale Doshi-Velez¹, Leo Anthony Celi^{3,6,7}

4

* These authors contributed equally to the preparation of the manuscript

¹ Paulson School of Engineering and Applied Sciences, Harvard University

² Institute for Medical Engineering and Science, MIT

³ Laboratory for Computational Physiology, Harvard-MIT Health Sciences & Technology, MIT

⁴ Department of Surgery and Cancer, Faculty of Medicine, Imperial College London

⁵ Department of Bioengineering, Imperial College London

⁶ Division of Pulmonary, Critical Care and Sleep Medicine, Beth Israel Deaconess Medical Center

⁷ MIT Critical Data

5 **Abstract**

6 In this Comment, we provide guidelines for reinforcement learning for patient treatment
7 decisions that we hope will accelerate the rate at which observational cohorts can inform
8 healthcare practice in a safe, risk-conscious manner.

9 From sepsis warning systems to identifying subtle disease signals in medical images, artificial
10 intelligence (AI) is poised to transform healthcare for the better¹. However, AI is not a panacea,
11 and if used improperly, these systems can replicate our bad practices rather than improve them.

12
13 Reinforcement Learning (RL) is a subfield of AI that provides tools to optimize sequences of
14 decisions for long-term outcomes. For example, faced with a patient with sepsis, the intensivist
15 must decide if and when to initiate and adjust treatments such as antibiotics, intravenous fluids,
16 vasopressor agents, and mechanical ventilation. Each choice affects the patient's survival at the
17 end of the hospital stay, quality of life upon recovery, and so on. While the RL approaches used
18 to optimize treatment sequences vary, they all fall into a common framework. RL algorithms
19 take as input sequences of interactions (called histories) between the decision-maker and their
20 environment. At every decision-point, the RL algorithm chooses an action according to its policy
21 and receives new observations and immediate outcomes (often called rewards).

22
23 In the context of healthcare, RL has been applied to optimizing anti-retroviral therapy in HIV²,
24 tailoring anti-epilepsy drugs for seizure control³, and determining the best approach to managing
25 sepsis⁴. In contrast to more common uses of AI such as one-time predictions, the output
26 (decisions) of an RL system affects both the patient's future health and future treatment
27 options⁵. As a result, long-term effects are harder to estimate (Figure 1a).

28
29 To illustrate the potential pitfalls in reinforcement learning, we use the example of sepsis
30 management, for which there remains wide variability in the way clinicians make decisions. In
31 the context of sepsis, a history may include a patient's vital signs and laboratory tests. The
32 actions are all the treatments available to the clinician, including medications and interventions.
33 The rewards require clinician input: they should represent the achievement of desirable tasks
34 such as stabilization of vital signs or survival at the end of the stay. By weighing different

35 rewards, an RL-algorithm could be designed to target short-term outcomes, such as liberation
36 from mechanical ventilation, or longer-term outcomes, such as prevention of permanent organ
37 damage. Note that defining short-term goals is not straightforward, since ideal sepsis
38 resuscitation targets remain elusive⁶.

39

40 We discuss three key questions that should be considered when reading an RL study. These
41 questions uncover limitations when making quantitative performance claims about RL-learned
42 algorithms from observational data.

43

44 **Is the AI given access to all the variables that influence decision making?**

45 A clinician could not be expected to make good decisions about a patient's vasopressor
46 medication dosing without knowing about the patient's co-morbid cardiac condition as well as
47 what has transpired in the last 24 hours, and neither can an AI. To estimate the quality of a new
48 treatment policy based on historical data, it is vital to take into account any information that was
49 used by clinicians in their decision making—failing to do so may result in estimates that are
50 confounded by spurious correlation. For example, severely sick septic patients may receive
51 fluids earlier, yet have worse outcomes than healthier patients which is clearly a result of them
52 being sicker in the first place. This difference in outcomes may lead an analysis that associates
53 earlier fluid administration with worse outcomes if not properly adjusted for clinical context.

54 Adjusting for confounding is challenging when validating the average treatment effect of a single
55 decision⁷; this problem becomes significantly harder when decisions are made in sequence. It is
56 thus important to be conscientious of possible confounding factors when reading an RL study
57 even more so than for standard prediction studies, as the sequential nature of the problem could
58 lead to confounding effects on the long as well as the short term.

59

60 **Effective cohort size: How big was that big data, really?**

61 When evaluating the quality of an RL algorithm retrospectively, the choice of the proposed
62 treatment policy affects the effective sample size. This occurs because most approaches for
63 evaluating RL policies from observational data weigh each patient's history based on whether
64 the clinician decisions match the decisions of the policy proposed by the RL algorithm⁸. The
65 reliability (variance) of the treatment quality estimate depends on the number of patient histories
66 for which the proposed and observed treatment policies agree — a quantity known as the
67 effective sample size. The possibilities for mismatch between the actual decision and the
68 proposed decision grow with the number of decisions in the patient's history, and thus RL-based
69 evaluations are especially prone to having small *effective* sample sizes (Figure 1b).

70

71 For example, we found that the effective sample size for a sepsis management policy on a
72 cohort of 3855 patients was only a few dozens⁹. In general, the effective sample size will be
73 larger if the learned policies are close to the clinician policies, suggesting that RL with
74 observational data will be most reliable for refining existing practices rather than discovering
75 new treatment approaches.

76

77 **Will the AI behave prospectively as intended?**

78 Even if the AI has access to all the important variables and the evaluation was perfect, errors in
79 problem formulation or data processing can lead to poor decisions. Simplistic reward functions
80 may neglect long-term effects for meaningless gains: for example, rewarding only blood
81 pressure targets may result in an AI that causes long-term harm by excessive dosing of
82 vasopressors. Errors in data recording or preprocessing may introduce errors in the reward
83 signal, misleading the RL algorithm. Finally, the learned policy may not work well at a different
84 hospital or even in the same hospital a year later if treatment standards shift.

85

86 Thus, it is essential to interrogate RL-learned policies to assess whether they will behave

87 prospectively as intended. An increasing body of work on interpretable machine learning
88 enables such introspection¹⁰.

89

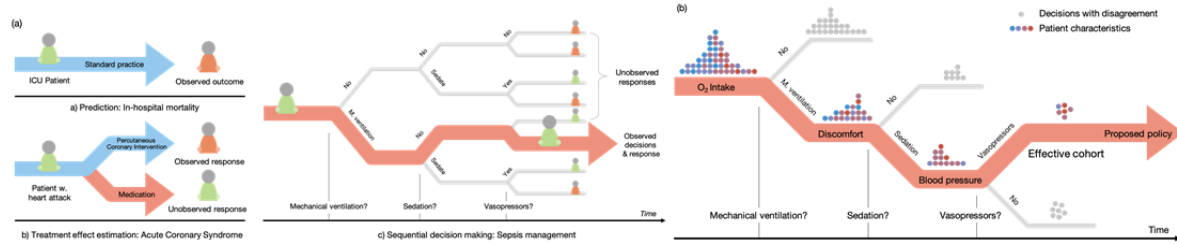
90 **Toward standard practice**

91 Together, big data and RL provide unique opportunities for optimizing treatments in healthcare,
92 especially those undertaken in sequence. However, to realize this potential, we must exercise
93 caution and due diligence in their application and evaluation.

94

95

96 **Figures**



97 **Figure 1. (a)** Prediction, treatment effect estimation and sequential decision-making tasks.
98 These tasks are progressively harder to solve based on observational data. In classical
99 prediction tasks, only a single outcome for a patient is considered—the result of following
100 standard practice without interventions from the analyst. Here, we use the common example of
101 predicting 48h in-hospital mortality. In treatment effect estimation, we must also reason about
102 what would happen under alternative unobserved interventions. Consider for example choosing
103 between performing catheterization on a patient with cardiac arrest, or placing them on
104 medication. To perform sequential decision making, such as for sepsis management, treatment
105 effect estimation must be solved at a much grander scale—every possible combination of
106 interventions could be considered to find an optimal treatment policy. **(b)** Effective sample size
107 in off-policy evaluation. Each dot represents a single patient at each stage of treatment, its color
108 indicating the patient’s characteristics. The more decisions are performed in sequence, the
109 likelier it is that a new policy disagrees with the one used to learn from. We illustrate
110 disagreement by grayed out decision points. Using only samples for which the old policy agrees
111 with the new results in a small effective sample size and a biased cohort, as illustrated by the
112 difference in color distribution in the original and final cohort.
113
114

115 1 Obermeyer, Z. & Emanuel, E. J. Predicting the future—big data, machine learning, and
116 clinical medicine. *The New England journal of medicine* **375**, 1216 (2016).

117 2 Parbhoo, S., Bogojeska, J., Zazzi, M., Roth, V. & Doshi-Velez, F. Combining Kernel and
118 Model Based Learning for HIV Therapy Selection. *AMIA Summits on Translational
119 Science Proceedings* **2017**, 239 (2017).

120 3 Guez, A., Vincent, R. D., Avoli, M. & Pineau, J. in *AAAI* 1671-1678 (2008).

121 4 Komorowski, M., Celi, L. A., Badawi, O., Faisal, A. & Gordon, A. The intensive care AI
122 clinician learns optimal treatment strategies for sepsis. *Nature Medicine* (2018 (In press)).

123 5 Chakraborty, B. & Moodie, E. *Statistical methods for dynamic treatment regimes*.
124 (Springer, 2013).

125 6 Simpson, N., Lamontagne, F. & Shankar-Hari, M. Septic shock resuscitation in the first
126 hour. *Current opinion in critical care* **23**, 561-566 (2017).

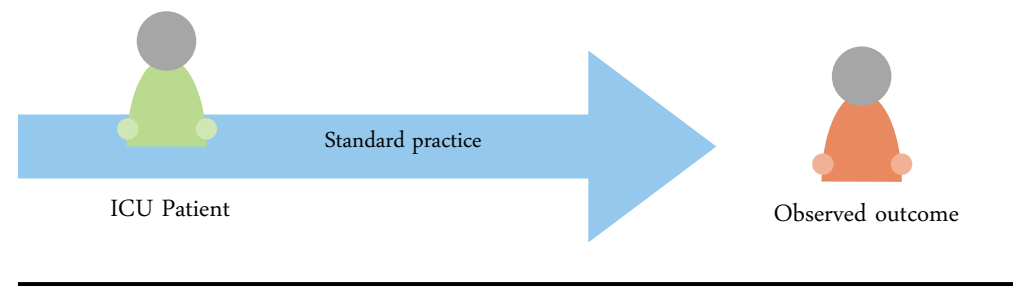
127 7 Johansson, F., Shalit, U. & Sontag, D. in *International Conference on Machine Learning*
128 3020-3029 (2016).

129 8 Precup, D., Sutton, R. S. & Singh, S. P. in *International Conference on Machine
130 Learning* 759-766 (2000).

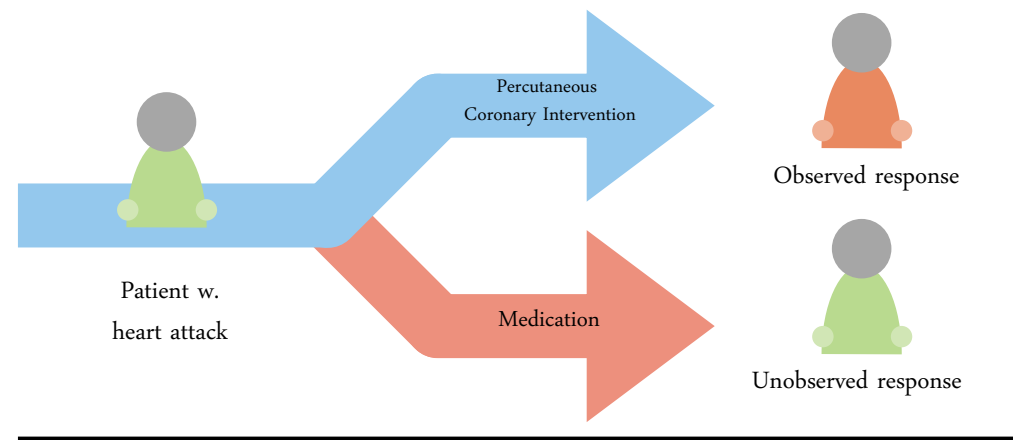
131 9 Gottesman, O. *et al.* Evaluating Reinforcement Learning Algorithms in Observational
132 Health Settings. *arXiv preprint arXiv:1805.12298* (2018).

133 10 Doshi-Velez, F. & Kim, B. Towards a rigorous science of interpretable machine learning.
134 *arXiv preprint arXiv:1702.08608* (2017).
135

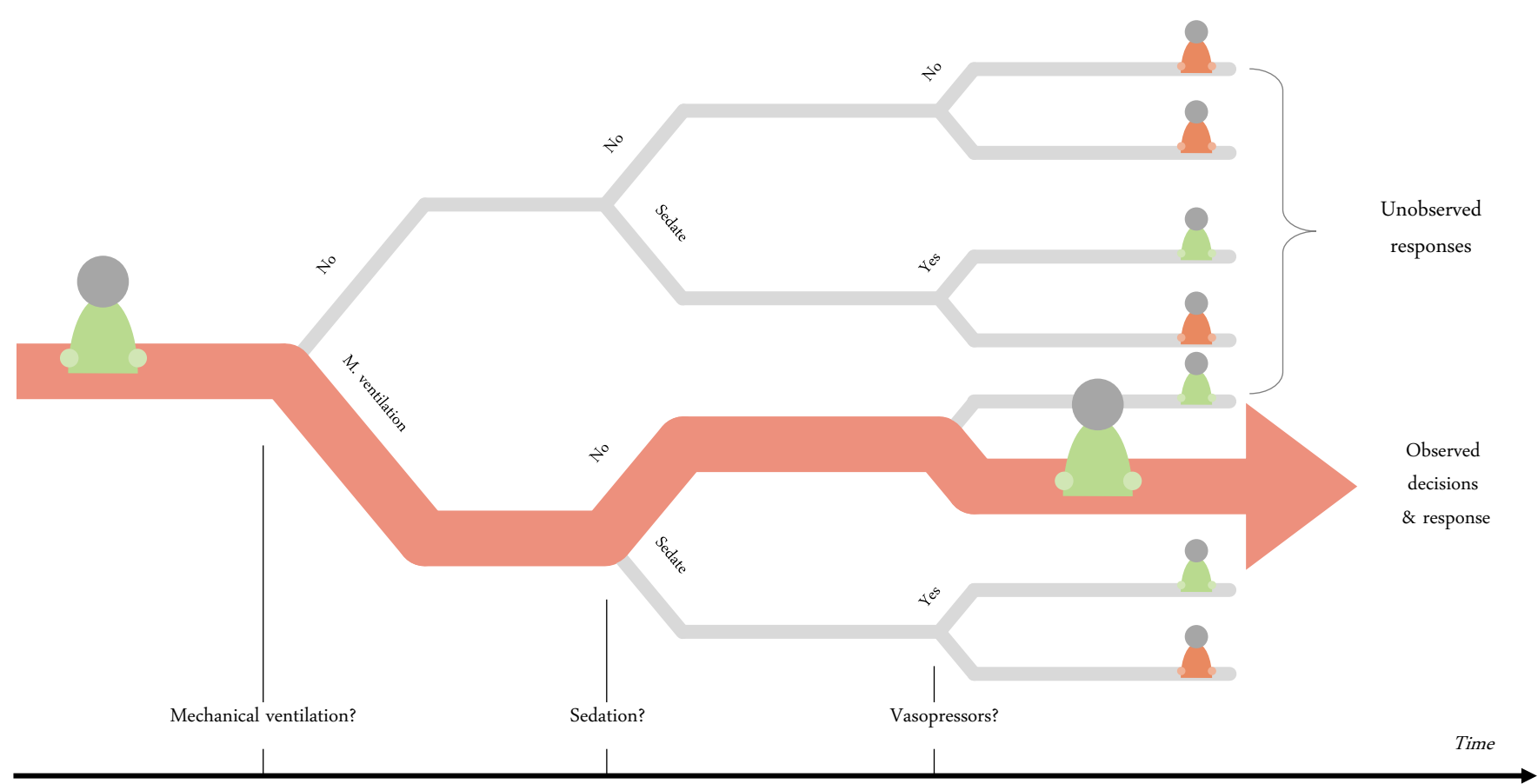
(a)



a) Prediction: In-hospital mortality



b) Treatment effect estimation: Acute Coronary Syndrome



c) Sequential decision making: Sepsis management

(b)

● Decisions with disagreement
●●● Patient characteristics

