# Challenges with reinforcement learning model transportability for sepsis treatment in emergency care

Check for updates

Peter C. Nauka[1], Jason N. Kennedy[2,3], Emily B. Brant[2], Matthieu Komorowski[4], Romain Pirracchio[5], Derek C. Angus[2,3] & Christopher W. Seymour[2,3,6] ✉

Pivotal moments in sepsis care occur in the emergency department (ED), however, and it is unclear whether ED data is adequate to inform reinforcement learning (RL) models. We evaluated the early opportunity for the AI Clinician, a validated ICU-based RL-model, as a use case. Amongst emergency sepsis patients, model parameters were often missing and invariably measured. Current iterations of RL-models trained on ICU data face challenges in emergency sepsis care.

Sepsis is deadly and accounts for nearly one in five hospital deaths[1–5]. It is a heterogenous syndrome of infection and life-threatening organ dysfunction for which the current "one size fits all" treatment is inadequate[3,6]. Development of intelligent tools to provide tailored treatment recommendations for sepsis resuscitation are needed.

Artificial intelligence (AI) is a promising tool to assist clinical decisions about sepsis treatment[7]. Reinforcement learning (RL) models such as the "AI Clinician" are proposed to guide hemodynamic support[8]. The AI Clinician uses a partially imputed set of 48 clinical features including vital signs, demographics, laboratory parameters, medications, and fluid balance to recommend an precision treatment policy of vasopressors and intravenous fluids every four hours. Following the AI Clinician's recommendations was associated with improved patient outcomes and decreased mortality of septic patients during intensive care. However, the most essential diagnostic and treatment decisions for sepsis patients are made in the emergency department (ED). To date, few RL models target care in the emergency department due to challenges of incomplete data which is invariably measured and unclear diagnosis of sepsis.

In a prospective cohort study of sepsis patients, we characterized the challenges for RL tools that were trained on intensive care data, using the AI Clinician as a use case, across four domains, (i) feature missingness, (ii) measurement cadence, (iii) diagnostic uncertainty, and (iv) treatment variability in the first six hours after arrival.

Among all adults presenting to the ED, 71,272 (mean age 68 years (SD, 17 years), 48% male) met Sepsis-3 criteria within six hours of hospital presentation (Table 1). The median Elixhauser comorbidity score was 1.0 [IQR: 0.0, 2.0], and median sequential organ failure assessment (SOFA) score was 3 [IQR: 2, 5]. Less than half (36%, $n = 25,579$) required intensive care admission, received mechanical ventilation (10%, $n = 7377$), or died at 90 days (20%, $n = 14,082$).

To assess feature missingness for the AI Clinician use case, we evaluated 36 of 44 available, non-calculated features mapped in one-hour time epochs. Fixed variables such as age, gender, weight, and Elixhauser score were non-missing (Fig. 1, Supplementary Table 1). Of the seven vital signs (Glasgow Coma Scale score, heart rate, systolic blood pressure, diastolic blood pressure, respiratory rate, temperature, and arterial partial pressure of oxygen), missingness ranged from 0% to 79% (Supplementary Table 2). For instance, the heart rate was measured in the first hour in 2% (1401) of patients, repeated in 70% ($n = 50,193$), and missing in 0.5% ($n = 365$).

Of the 19 laboratory parameters, feature missingness ranged from 13% to 85%. Laboratory variables such as carbon dioxide ($CO_2$) (84%, $n = 59,960$), arterial partial pressure of oxygen ($PaO_2$) (85%, $n = 60,711$), and pH (84%, $n = 59,853$) were commonly missing, and of those measured in the first hour, repeat measurement ranged from 1% to 13% (Supplementary Table 2). For instance, glucose was missing in 14% ($n = 9932$), measured in the first hour only in 41% ($n = 29,085$) of patients, and repeated in 13% of patients ($n = 8980$). Organ support variables were routinely available. However, the level of support was often missing in the first hour, including $FiO_2$ (present in 2% of patients, $n = 1086$). Fewer than 1% ($n = 36$) of patients had complete measurement of all features by hour six.

The diagnosis of sepsis was often uncertain. For example, only 26% ($n = 18,752$) and 45% ($n = 32,209$) of patients met Sepsis-3 criteria in the first and second hour, respectively. In the first hour, 46% of patients ($n = 32,604$) had a SOFA score ≥ 2, 44% ($n = 31,137$) had body fluid culture obtained, and 4% ($n = 2718$) were administered antibiotics. These cumulative rates increased throughout the six-hour measurement window (Fig. 2).

[1]Division of Pulmonary, Allergy, Critical Care, and Sleep Medicine, Department of Medicine, University of Pittsburgh Medical Center, Pittsburgh, PA, USA. [2]Department of Critical Care Medicine, University of Pittsburgh School of Medicine, Pittsburgh, PA, USA. [3]Clinical Research, Investigation and Systems Modeling of Acute Illness (CRISMA) Center, Pittsburgh, PA, USA. [4]Department of Surgery and Cancer, Imperial College London, London, UK. [5]Department of Anesthesia and Perioperative Medicine, Zuckerberg San Francisco General Hospital and Trauma Center, University of California San Francisco, San Francisco, CA, USA. [6]Department of Emergency Medicine, University of Pittsburgh School of Medicine, Pittsburgh, PA, USA. ✉e-mail: seymourcw@upmc.edu

**Table 1 | Clinical characteristics of cohort**

| Variable | Cohort |
|---|---|
| No. of patients | 71,272 |
| Age [yrs]: mean (SD) | 68 (17) |
| Sex[a] [male]: no. (%) | 34,384 (48%) |
| Elixhauser comorbidity[b]: median [IQR] | 1 [0, 2] |
| Maximum 6-h SOFA score: median [IQR] | 3 [2, 5] |
| Race[c]: no. (%) | |
| White | 59,245 (85%) |
| Black | 8304 (12%) |
| Other | 2441 (3%) |
| Hospital use and outcomes | |
| Hospital length of stay [days]: median [IQR] | 6 [4, 10] |
| Mechanical ventilation, no. (%) | 7377 (10%) |
| Glucocorticoids during hospitalization, no. (%) | 8957 (13%) |
| Surgery during hospital stay, no. (%) | 15,071 (21%) |
| ICU admission, no. (%) | 25,579 (36%) |
| Mortality, 90-day, no. (%) | 14,082 (20%) |

*ICU* intensive care unit, *no.* number, *SD* standard deviation, *IQR* interquartile range, *yrs* years, *SOFA* sequential organ failure assessment.

[a]A total of four patients had missing sex parameter.

[b]Elixhauser is a method of categorizing comorbidities of patients based on the International Classification of Diseases (ICD) diagnosis codes found in administrative data, ranging from 0 to 31 (*n* = 151 missing).

[c]A total of 1282 patients had missing race data. Percentages correspond to total available. Other race corresponds to Chinese, Filipino, Hawaiian, American Indian/Alaskan, Asian, Hawaiian/Other Pacific Islander, Middle Eastern, Native American, Not specified, or Pacific Islander.

Treatment variability was commonly present. Despite meeting Sepsis-3 criteria, almost half of patients (49%, *n* = 35,222) did not receive any intravenous fluids or vasopressors in the six-hour period (Supplementary Fig. 1, Supplementary Table 3). Of those treated, 46% (*n* = 32,835) received only intravenous fluids, 1% (*n* = 776) received only vasopressors, and 3% (*n* = 2439) were documented to receive both. The frequency and amount of treatment increased from hour one to six (Supplementary Fig. 2). For example, clinicians administered intravenous fluids in 17% (*n* = 12,172) of patients during hour one, 30% (*n* = 21,039) in hour two, and 50% (*n* = 35,274) by hour six (Supplementary Figs. 2, and 3, Supplementary Tables 4, and 5). Meanwhile, amongst patients receiving vasopressors, mean norepinephrine equivalent dose increased from 0.36 μg/kg/min (95%CI: 0.32–0.40) to 0.46 μg/kg/min (95%CI: 0.44–0.49) by hour six. (Supplementary Fig. 3, Supplementary Tables 4, and 5).

These findings highlight several challenges in the development and implementation of RL-based clinical decision support during emergency care. Critical conditions like sepsis require time-sensitive treatment which are crucial for improving outcomes. Despite the potential of current RL models such as the Clinician AI, our study reveals significant limitations in feature availability, inconsistent laboratory measurement schedules, and delayed diagnostic certainty. Overcoming these challenges will be essential for the successful application of RL methods in emergency care.

The current approach to train clinical decision support models leans heavily on retrospective data in four to six hour epochs[8,9]. Our data suggest current versions of learning policies using single value, carry and hold imputation is not ideal for one-hour epochs in emergency care. It is likely that sampling frequency is strongly associated with non-random missingness[10]. Given the limits of current imputation methods, automated systems that transcribe vital signs or standardize clinical data exchange via tools such as Health Level Seven International (HL7) Fast Healthcare Interoperability Resources (FHIR) could be a solution. Although these standards are evolving, the mapping of features to FHIR could improve the internal validity of the data scaffold on which RL models are

trained—both between locations within a hospital as well as between hospitals themselves. Another approach to improve RL model performance is incorporating uncertainty within RL decisional framework. For example, a partially observable Markov decision process modelling (POMDP) is a method that uses available information to make decisions under uncertainty[11]. This approach could facilitate timely policy recommendations even in emergency care where key clinical features and sepsis diagnosis are obtuse. Beyond improving data granularity or incorporating uncertainty, the selection of a cohort of emergency sepsis patients eventually admitted to the ICU may lead to RL models less hampered by data paucity and diagnostic uncertainty.

The transportability of RL models trained on ICU data (like the Clinician AI) to emergency care of sepsis is concerning. First, model performance varies by healthcare setting and case mix. The distribution of treatments in alternative settings from the ICU, perhaps like the administration (or not) of vasopressors, pose model transportability issues. There is the potential for unexpected outcomes if RL models are, thus, applied to "off label" in populations[12,13]. Second, the clinical data which underpin RL models can be biased and transporting those assumptions to new populations may be detrimental[14–17]. Third, sepsis is a challenging diagnosis with several axes of complexity (e.g., unknown initial time zero, chronic organ dysfunction) that stress the direct transportability of RL models trained in the ICU. New approaches to clinical validation (rather than traditional randomized trials) may be required for RL technology to be effectively and safely deployed at emergency departments. Rigorous post-development and post-deployment monitoring will be required to ensure that RL-based models can improve patient outcomes.

This study has several limitations. First, we used a single US health care system, which reduces generalizability of our findings in low-middle income or alternative settings. Second, the data originated from 2010 to 2014, and the patterns of vasopressor and intravenous fluid may have changed as clinical practice guidelines evolve. Third, we assumed a lack of documentation for certain respiratory and hemodynamic support variables indicated an absence of measurement, which is a common but unverified assumption[18–20]. Fourth, even in a well-researched electronic health record dataset, we were unable to map the full set of variables from the AI Clinician highlighting the need for future RL model parsimony[20].

In conclusion, the current approach to transport reinforcement learning models for sepsis care from the ICU to the emergency department is inadequate. Common factors such as missing data, diagnostic uncertainty, and invariably sampled parameters suggest new models are required.

## Methods

To broadly understand the opportunities and challenges for AI clinical decision support during emergency care of sepsis, we specifically chose the AI Clinician as a use case[8]. The AI Clinician is a reinforcement learning model using 46 clinical variables validated in over 96,000 patients. The study was approved by the University of Pittsburgh Institutional Review Board (IRB 20010238) with a waiver of informed consent due to the research posing only minimal risk.

### Patients

The cohort includes adult patients presenting to an integrated healthcare system of 14 academic and community hospitals in Western Pennsylvania from January 1, 2010 to December 31, 2014. Patients must have met Sepsis-3 criteria within six hours of arrival, defined as: i. Suspected infection, and ii. presence of organ dysfunction as measured by a SOFA score of two or more[20]. According to Sepsis-3, suspected infection is operationalized as the administration of antimicrobial therapy (oral or parenteral) and body fluid culture sampling (e.g., blood, urine or cerebrospinal fluid).

### Data

We extracted patient-level data from the electronic health records (CERNER, Inc), including demographics, Elixhauser comorbidity score, vital signs, laboratory measurements, microbiology,
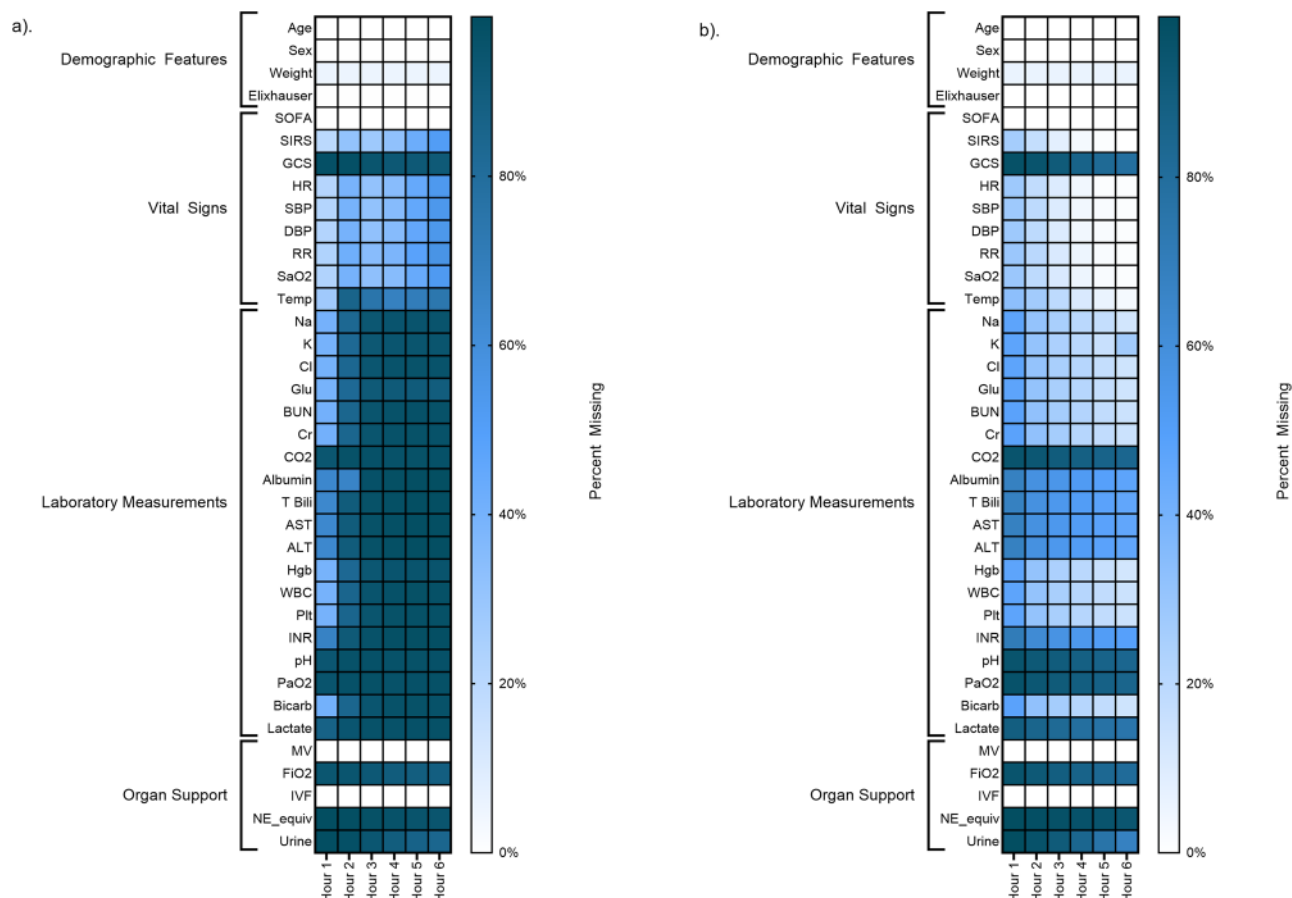
**Fig. 1 | Heatmap of missing feature by hour from admission.** Darker blue squares demonstrate a higher missing rate of individual AI Clinician feature. Panel **a** shows the percentage missing per hour. Panel **b** shows cumulatively missing rate by hour. Abbreviations: Elix Elixhauser Comorbidity Score, SOFA Sequential organ failure score, SIRS systemic inflammatory response score, GCS Glasgow-Coma scale, HR heart rate, SBP systolic blood pressure, DBP diastolic blood pressure, RR respiratory rate, SaO2 peripheral oxygen saturation, Temp temperature, Na sodium, K potassium, Cl chloride, Glu glucose, BUN blood urea nitrogen, Cr creatinine, CO2 carbon dioxide, Alb albumin, T Bili total bilirubin, AST asparagine transferase, ALT alanine transferase, Hgb hemoglobin, WBC white blood cell count, Plt platelet, INR international normalized ratio, PaO2 arterial partial oxygen, Bicarb bicarbonate, MV mechanical ventilation, FiO2 fractional inspired oxygen, IVF intravenous fluids, NE_equiv, norepinephrine equivalents.
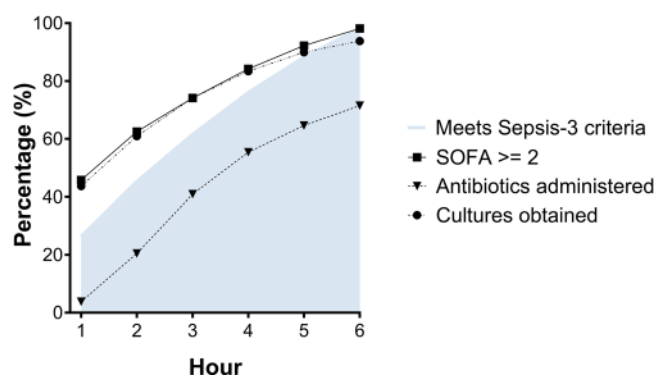


**Fig. 2 | Percentage of cohort meeting Sepsis-3 criteria over study period.** Cumulative percentage of cohort meeting Sepsis-3 criteria (light blue curve) during initial six hours after presentation to emergency department. Individual components (SOFA score ≥2, fluid culture sampling, antibiotic administration) demonstrated as dashed line. Percentage of patients meeting SOFA ≥2 and Sepsis-3 criteria does not equal 100% given averaging of variables repeated in multiple epochs. Abbreviations: SOFA Sequential organ failure score.

intravenous fluids, and vasopressor administration. The initial AI Clinician model contained a total of 48 variables, but 12 were unavailable. Eight variables (magnesium, calcium, ionized calcium, PT, PTT, PaCO2 and base excess, readmission to ICU) were not available in the data set while others (e.g., shock index, PaO2/FiO2 ratio, cumulative fluid balance, mean blood pressure) could be theoretically computed but are reported as missing. Race was determined from initial patient registration data and categorized according to the Centers for Medicare & Medicaid Services EHR meaningful use data set. Comorbidities used for the Elixhauser comorbidity score were classified by ICD-9 and ICD-10 hospitalization diagnosis codes[21].

Data was organized in non-overlapping, one-hour blocks from initial presentation to the ED. For variables with multiple measurements in the one-hour block, we used the mean value of all measurements. Intravenous fluid volumes were measured each hour (mL), and vasopressor administration was converted to norepinephrine-equivalents[22].

The SOFA score was calculated for each one-hour time block using a standard approach with missing parameters assumed to be normal. Patient outcomes included in hospital and 90-day mortality, determined by inpatient discharge disposition and monthly Social Security Index.

## Analysis

To understand the opportunities and challenges for the AI Clinician and similar models, we analyzed four domains, (i) feature missing, (ii) measurement cadence, (iii) diagnostic uncertainty and (iv) treatment variability. Measurement cadence was summarized by reporting the percentage of each variable which was never measured, measured in the first hour of clinical care only, those in later hours only, and those repeated during later epochs. Diagnostic uncertainty was characterized by the cumulative proportion of patients meeting Sepsis-3 criteria in each hour epoch. The individual components of Sepsis-3 (e.g., SOFA > 2, antimicrobial administration, and body fluid sampling) were reported. We summarized AI Clinician treatment variation (vasopressors and intravenous fluid administration) by proportion of patients receiving each treatment as well as cumulative hourly dose. Categorical data was presented as either N or %, as appropriate. Continuous data was summarized as mean (SD) or median [IQR]. We illustrated feature missingness using a heatmap across one-hour epochs. GraphPad Prism version 10 (GraphPad Software) was used, and analyses performed with Stata version 18.0 (StataCorp LLC).

## Data availability

The dataset is available upon request from the corresponding authors and subsequent institutional board approval.

## Code availability

The code is available upon request from corresponding authors.

## References

1. Rudd, K. E. et al. Global, regional, and national sepsis incidence and mortality, 1990-2017: analysis for the Global Burden of Disease Study. *Lancet* **395**, 200–211 (2020).
2. Evans, L. et al. Executive summary: surviving sepsis campaign: international guidelines for the management of sepsis and septic shock 2021. *Crit. Care Med.* **49**, 1974–1982 (2021).
3. Seymour, C. W. et al. Time to treatment and mortality during mandated emergency care for sepsis. *N. Engl. J. Med.* **376**, 2235–2244 (2017).
4. Ferrer, R. et al. Empiric antibiotic treatment reduces mortality in severe sepsis and septic shock from the first hour: results from a guideline-based performance improvement program. *Crit. Care Med.* **42**, 1749–1755 (2014).
5. Rhee, C. et al. Incidence and trends of sepsis in US Hospitals using clinical vs claims data, 2009-2014. *JAMA* **318**, 1241–1249 (2017).
6. Barbash, I. J. et al. Treatment patterns and clinical outcomes after the introduction of the medicare sepsis performance measure (SEP-1). *Ann. Intern. Med.* **174**, 927–935 (2021).
7. Boussina, A. et al. Impact of a deep learning sepsis prediction model on quality of care and survival. *npj Digital Med.* **7**, 14 (2024).
8. Komorowski, M., Celi, L. A., Badawi, O., Gordon, A. C. & Faisal, A. A. The Artificial Intelligence Clinician learns optimal treatment strategies for sepsis in intensive care. *Nat. Med.* **24**, 1716–1720 (2018).
9. Peine, A. et al. Development and validation of a reinforcement learning algorithm to dynamically optimize mechanical ventilation in critical care. *npj Digital Med.* **4**, 32 (2021).
10. Laudermilch, D. J., Schiff, M. A., Nathens, A. B. & Rosengart, M. R. Lack of emergency medical services documentation is associated with poor patient outcomes: a validation of audit filters for prehospital trauma care. *J. Am. Coll. Surg.* **210**, 220–227 (2010).
11. Chadès, I., Pascal, L. V., Nicol, S., Fletcher, C. S. & Ferrer-Mestres, J. A primer on partially observable Markov decision processes (POMDPs). *Methods Ecol. Evol.* **12**, 2058–2072 (2021).
12. Classen, D. C., Longhurst, C. & Thomas, E. J. Bending the patient safety curve: how much can AI help? *npj Digital Med.* **6**, 2 (2023).
13. Wong, A. et al. External validation of a widely implemented proprietary sepsis prediction model in hospitalized patients. *JAMA Intern. Med.* **181**, 1065–1070 (2021).
14. Celi, L. A. et al. Sources of bias in artificial intelligence that perpetuate healthcare disparities-A global review. *PLoS Digit Health* **1**, e0000022 (2022).
15. Ranard, B. L. et al. Minimizing bias when using artificial intelligence in critical care medicine. *J. Crit. Care* **82**, 154796 (2024).
16. Zack, T, et al. Assessing the× potential of GPT-4 to perpetuate racial and gender biases in health care: a model evaluation study. *Lancet Digit Health* **6**, e12–e22 (2024).
17. Iqbal, U., Hsu, Y.-H. E., Celi, L. A. & Li, Y.-C. J. Artificial intelligence in healthcare: opportunities come with landmines. *BMJ Health Care Inform* **31**, e101086 (2024).
18. Knaus, W. A., Draper, E. A., Wagner, D. P. & Zimmerman, J. E. APACHE II: a severity of disease classification system. *Crit. Care Med.* **13**, 818–829 (1985).
19. Smith, G. B., Prytherch, D. R., Meredith, P., Schmidt, P. E. & Featherstone, P. I. The ability of the National Early Warning Score (NEWS) to discriminate patients at risk of early cardiac arrest, unanticipated intensive care unit admission, and death. *Resuscitation* **84**, 465–470 (2013).
20. Seymour, C. W. et al. Assessment of clinical criteria for sepsis: for the Third International Consensus Definitions for Sepsis and Septic Shock (Sepsis-3). *JAMA* **315**, 762–774 (2016).
21. Elixhauser, A., Steiner, C., Harris, D. R. & Coffey, R. M. Comorbidity measures for use with administrative data. *Med. Care* **36**, 8–27 (1998).
22. Brown, S. M. et al. Survival after shock requiring high-dose vasopressor therapy. *Chest* **143**, 664–671 (2013).

## Author contributions

Christopher Seymour had full access to all the data in the study and takes responsibility for the integrity of the data and the accuracy of the data analysis. All authors (P.C.N., J.N.K., E.B.B., M.K., R.P., D.C.A., C.W.S.) have read and approved this manuscript. Concept and design: Nauka, Brandt, Seymour Acquisition, analysis, or interpretation of data: Nauka, Brandt, Kennedy, Seymour Drafting of the manuscript: Nauka, Seymour Critical revision of the manuscript for important intellectual content: All authors. Statistical analysis: Nauka, Kennedy Administrative, technical, or material support: Seymour Supervision: Seymour.

## Competing interests

Dr. Seymour reports grants from the NIH and personal fees from Inotrem, Deepull, Beckman Coulter, and Octapharma outside the submitted work. Dr. Komorowski would like to disclose consulting fees from Philips Healthcare and speaker honoraria from GE Healthcare. The remaining authors (Dr Nauka, Mr Kennedy, Dr Brandt, Dr Pirrachio, Dr Angus) declare no competing interests.

## Additional information

**Supplementary information** The online version contains supplementary material available at https://doi.org/10.1038/s41746-025-01485-6.

**Correspondence** and requests for materials should be addressed to Christopher W. Seymour.

**Reprints and permissions information** is available at http://www.nature.com/reprints