

Clinical Deployment Environments: The Five Pillars of Translational Machine Learning for Health

Steve Harris^{1,2*}† Tim Bonnici^{1,2†} Thomas Keen¹ Watjana Lilaonitkul¹ Mark J White⁴ Nel Swanepoel³

¹*Institute of Health Informatics, University College London, London, UK*

²*Department of Critical Care, University College London Hospital, London, UK*

³*Centre for Advanced Research Computing, University College London, London, UK*

⁴*Digital Healthcare, University College London Hospital, London, UK*

Correspondence*:

Corresponding Author

steve.harris@ucl.ac.uk

†These authors have contributed equally to this work and share first authorship

3

4 ABSTRACT

5 Machine Learning for Health (ML4H) has demonstrated efficacy in computer imaging and other
6 self-contained digital workflows, but has failed to substantially impact routine clinical care. This is
7 no longer because of poor adoption of Electronic Health Records Systems (EHRS), but because
8 ML4H needs an infrastructure for development, deployment and evaluation within the healthcare
9 institution. In this paper, we propose a design pattern called a Clinical Deployment Environment
10 (CDE). We sketch the five pillars of the CDE: (1) real world development supported by live data
11 where ML4H teams can iteratively build and test at the bedside (2) an ML-Ops platform that brings
12 the rigour and standards of continuous deployment to ML4H (3) design and supervision by those
13 with expertise in AI safety (4) the methods of implementation science that enable the algorithmic
14 insights to influence the behaviour of clinicians and patients and (5) continuous evaluation that
15 uses randomisation to avoid bias but in an agile manner. The CDE is intended to answer the same
16 requirements that bio-medicine articulated in establishing the translational medicine domain. It
17 envisions a transition from 'real-world' data to 'real-world' development.

18 **Keywords:** translational medicine, machine learning, health informatics, ml-ops, safety, artificial intelligence

INTRODUCTION

19 Bold claims and huge investments suggest Machine Learning (ML) will transform healthcare.(1) High
20 impact publications showcase precision models that predict sepsis, shock, and acute kidney injury.(2, 3, 4)
21 Outside healthcare, tech titans such as AirBnB, Facebook, and Uber create value from ML despite owning
22 'no property, no content and no cars'.(5) Inspired by this, and very much aware of the flaws and unwarranted
23 variation in human decision making(6), government and industry are now laying heavy bets on ML for
24 Health (ML4H).(7, 8)

Widespread adoption of electronic health records (EHR) might be thought a sufficient prerequisite for this ambition. Yet while EHR adoption is growing at pace(9), those ML4H models that have reached the market rarely use the EHR. They are instead embedded in isolated digital workflows (typically radiology) or medical devices.(10) Here the context of deployment is static and self-contained (imaging), or fully specified (devices), and translation has proved easier to navigate.

In contrast, the EHR is in constant flux. Both the data and the data model are updating. New wards open, staffing patterns are adjusted and from time to time major incidents (even global pandemics) disrupt everything. There are multiple interacting users, and eventually there will be multiple interacting algorithms, and organizations will face the ML equivalent of poly-pharmacy.(11) Algorithms will require stewards.(12) Whilst the aforementioned high impact prediction models are developed on real-world data, this is not the same as real-world development. Data are either anonymized and analyzed offline, or moved out of the healthcare environment into an isolated Data Safe Haven (DSH) [also known as Trusted Research Environment (TRE)].(13) This separation is the first fracture leading to the oft-cited AI chasm(14) leaving the algorithms stranded on the laboratory bench.

A future that sees ML4H generate value from the EHR requires an alternative design pattern. TREs excel at meeting the needs of population health scientists but they do not have the full complement of features required to take an ML4H algorithm from bench-to-bedside. Using drug development as an analogy, a TRE is custom made for drug discovery not translational medicine.(15)

In this paper, we describe the functional requirements for a Clinical Deployment Environment (CDE) for translational ML4H. These requirements map closely to the classical components of translational medicine, but differ in that algorithms will require ongoing stewardship even after a successful deployment. The CDE is an infrastructure that manages algorithms with the same regard that is given to medicines (pharmacy) and machines (medical physics). Moreover, the value of ML4H will not just be from externally developed blockbuster models, but will also derive from specific and local solutions. Our vision of a CDE therefore enables both *development* and *deployment*.

Our CDE is supported by five pillars:

1. Real World Development
2. ML-Ops for Health
3. Responsible AI in practice
4. Implementation science
5. Continuous evaluation

We describe these pillars below alongside figures and vignettes reporting early local experience in our journey building this infrastructure.

1 REAL WORLD DEVELOPMENT

Real-world data (RW-Data) means the use of observational data at scale augmented by linking across multiple data sources to generate insights simply not available from isolated controlled clinical trials.(16) The FDA uses data from tens of millions of patients in its Sentinel programme¹ to monitor drug safety, and

¹ <<https://www.sentinelinitiative.org>>

the OpenSafely² programme in the UK generated impactful insights into COVID-19 within the first few months of the global pandemic.(17)

Given the sensitive nature of health data, these initiatives depend on expanding investment into TREs. (18) TREs are an example of 'data-to-modeler' (DTM) designs where data flows from source (primary, secondary, social care and elsewhere) to a separate, secure landing zone. Here research teams write the code to link, clean and analyze the data. Derived insights eventually return to the bedside through clinical guidelines and policy. To date, DTM is also the dominant design pattern in ML4H but this approach is fundamentally flawed.

It is flawed because it imposes a separation between the modeller and the end-user. ML4H is not concerned with better guidelines or policy but with better operational and clinical decision making. This requires the practitioner to work alongside the end-user because excellent offline model performance provides no guarantee of bedside efficacy. Algorithms with inferior technical performance may even provide greater bedside utility.(19, 20) An inverted 'modeler-to-data' (MTD) paradigm was initially proposed to reduce privacy concerns (data are no longer copied and shared but analyzed in situ(21)), but we see important additional value in that it forces 'real-world development' (RW-Dev) and enables the end-user to work with the modeler in rapid-cycle build-test-learn loops. This first pillar of the CDE is the equivalent of an *internal* TRE *within* the healthcare institution.(21)

RW-Dev has four functional sub-requirements that distinguish it from a TRE. (1) Firstly, data updates must match the cadence of clinical decision making. For most inpatient and acute care pathways, decisions are in real-time (minutes or hours) at the bedside or in the clinic. (2) Secondly, development using live data must be sandboxed and so the clinical system responsible for care delivery is protected (3) Thirdly, privacy must be managed such that teams are able to develop end-user applications that inevitably display patient identifiable information (PII) alongside the model outputs: an anonymous prediction is of little use to a clinician. (4) Fourthly, attention must be paid to developer ergonomics. Where development and deployment steps are separated physically (the TRE paradigm) or functionally (via different languages and technologies), ownership is often split between two different teams. One team prepares the raw data and develops the model, and another prepares the live data and deploys the model. We argue instead that the same team should be able develop *and* deploy. This should accelerate iteration, reduce cost and increase quality.(22)

We illustrate this idea with a description of our local real-world development platform in Figure 1, and provide an extended description in the Electronic Supplementary Material.

2 ML-OPS (FOR HEALTH)

Hitherto in ML4H, the data and the algorithm have been the 'celebrity couple'. State-of-the-art models trained on RW-Data deliver high profile publications.(4, 3) But only a tiny handful (fewer than 10 studies in a recent high quality systematic review of nearly 2000 ML4H publications(23)), were prospectively implemented. The standard offline 'data-to-modeler' (DTM) paradigm described above incurs a significant but 'hidden technical debt' that includes configuration, data collection and verification, feature extraction, analysis and process tools, compute and storage resource management, serving infrastructure, and monitoring.(24) In fact, the code for the underlying ML model is estimated to be at most 5% of the total code with the other 95% as additional code to make the system work. 'Glue-code', 'pipeline jungles',

² <<https://www.opensafely.org>>

and 'dead experimental codepaths' are some of the anti-patterns that make the transition into production costly and hazardous.³

Agencies such as the FDA⁴, EMA⁵, and MHRA⁶ are working toward safety standards for AI and machine learning, but the majority of these efforts derive from medical devices regulation. Treating Software as a Medical Device (SaMD) is appropriate where the algorithms operate within a constant and predictable environment (e.g. code embedded within a cardiac pacemaker). But, as already argued, ML4H models working with the EHR are likely to find themselves operating in a significantly more complex landscape. This inconstant environment where algorithms themselves may only have temporary utility has parallels to the commercial environment exploited so successfully by the tech giants.

These companies have cultivated an approach to model deployment called 'ML-Ops'. This combines the practices of 'DevOps' (a portmanteau of Software Development plus IT operations)(22) that focuses on the quality and speed with which software moves from concept to production, with robust data engineering and machine learning. A typical ML-Ops system monitors raw input data, checks for distribution drift, provides a feature store to avoid train/serve skew and facilitate collaboration between teams, and maintains an auditable and monitored model repository.(26) We present a prototype implementation interacting with the EHRS in Figure 2 (called FlowEHR).

This constant adjustment of algorithms based on their continuously measured quality and performance needs a workforce as well as a technology stack. Just as the safe delivery of medicines to the bedside is the central activity of a hospital pharmacy team, the safe delivery of algorithms will require the development of similarly skilled and specialized practitioners, and we should expect to see clinical ML-Ops departments in the hospital of the future. Others have made similar proposals and labeled this as "algorithmic stewardship" or "AI-QI".(12, 27) Similarly, the FDA is now proposing 'automatic Algorithmic Change Protocols' (aACP) and proposals have been advanced to guard against gradual deterioration in prediction quality ("biocreep").(28, 29)

3 RESPONSIBLE AI IN PRACTICE

Pillars 1 and 2 should engender well designed and well engineered algorithms, but they do not protect against the unintentional harm that AI may induce. Algorithms can only learn from a digital representation of the world that representation in turn cannot encode moral or ethical standards. Unfair outcomes, discrimination against sub-populations and bias are all reported shortcomings.(30) In a dynamic setting, risk can also arise in the form of degraded predictive performance over time. Models that modify clinician's behavior alter patient profiles by design, but predictive success today inevitably erodes future performance by rendering obsolete the historical patterns that drove the performance of the original model.(31) Responsible AI in practice requires a systems approach that preempts and safe-guards against these potential risks to patients. We highlight three promising responses to components of this challenge that need to become part of the risk management approach for ML4H.

³ One infamous example from the financial services sector saw a firm lose \$170,000 per second (more than \$400m in 45 minutes) when an outdated piece of code leaked into production. The firm in question was fined a further \$12m for "inadequate safeguards" allowing "millions of erroneous orders".(25)

⁴ <<https://www.fda.gov>>

⁵ <<https://www.ema.europa.eu/en>>

⁶ <<https://www.gov.uk/government/organisations/medicines-and-healthcare-products-regulatory-agency>>

3.1 Model explainability

We argue that model explainability (Explainable Artificial Intelligence [XAI]) methods need to be prioritized to help systematize and coordinate the processes of model troubleshooting by developers, risk-management by service providers, and system-checks by auditors.(32, 33, 34, 35) Most AI models that operate as ‘black-box models’ are unsuitable for mission-critical domains, such as healthcare, because they pose risk scenarios where problems that occur can remain masked and therefore undetectable and unfixable. We acknowledge recent critiques(36, 37) of explainability methods that argue the methods cannot yet be relied on to provide a determinate answer as to whether an AI-recommendation is correct. However, these methods do highlight decision-relevant parts of AI representations, and offer promise in measuring and benchmarking interpretability(38, 39). They are particularly promising for risk management as they can be used to structure a systematic interrogation of the trade-off between interpretability, model accuracy and the risk of model misbehavior.

3.2 Model fail-safes

Prediction models that map patient data to medically meaningful classes are forced to predict without the option to flag users when the model is unsure of an answer. To address this problem, there is good evidence that methods such as Bayesian deep learning and various uncertainty estimates (40) can provide promising ways to detect and refer data samples with high probability of misprediction for human expert review.(41, 42, 43) These fail safes, or selective prediction approaches should be designed into support systems to preempt and mitigate model misbehavior.(44, 45, 46, 47, 29) Of note, the European Commission High-Level Expert Group on AI presented guidelines for trustworthy AI in April 2019 with such recommendations: for systems that continue to maintain human-agency via a human-in-the-loop oversight. This may even permit less interpretable models to operate when implemented in conjunction with an effective fail-safe system.

3.3 Dynamic model calibration

As discussed, models that influence the evolution of its own future input data are at risk of performance deterioration over time due to input data shifts (48). In such cases, continual learning via calibration drift detection and model recalibration (27, 49) provides a promising solution but remains a challenging paradigm in AI. Recalibration with non-stationary incremental data can lead to catastrophic forgetting when the new data negatively interferes with what the model has already learned (50), or a convergence where the model just predicts its own effect and thus should not be updated (31). On the other hand, models can propose poor decisions because of the inherent biases found within the original dataset. In this case, dynamic model recalibration is unlikely to be sufficient and larger model revisions may be required. Here Pillar 1 (RW-dev) with suitable audit and monitoring via Pillar 2 (ML-Ops) will be required to overcome what would otherwise be a learning process encumbered by regulatory barriers.(51)

4 IMPLEMENTATION SCIENCE

A well designed, safe, and responsible AI algorithm may still be ineffective if it does not reach a modifiable target on the clinical pathway.(19) Unlike medications, algorithms can only effect health by influencing the behavior of clinicians and patients. This translational obstacle parallels the second arm of translational medicine (T2): implementation science.(15) Behavior change, in most instances, will be via a modification of the choice architecture (passive)(52, 53) or via interruptive alerts (active) embedded in the EHR.(53) Effective implementation requires a multi-disciplinary approach including human-computer interaction, behavioral science, and qualitative analysis.(54)

We strongly argue that this task will be more difficult if done offline and in isolation. Pillar 1 crucially permits not just tuning of the technical performance of the algorithm but rapid build-test-learn cycles that directly involve the target user and the clinical pathway in question. This approach will reduce costs and improve impact, sometimes leading to trade-offs which might appear surprising to those developing away from the bedside.(20, 11) This efficiency will again depend on the problem space: where the algorithmic target depends on information arising from the EHR rather than an isolated device or image, and where the pathway involves multiple end-users, then successful implementation will be near impossible if done sequentially (development then deployment) rather than iteratively.(54, 55) Academic health science centres must become design 'laboratories' where rapid prototyping at the bedside crafts the deployment pathway for *effectiveness* (T2) rather than just efficacy (T1).(15)

Investigations to define how system can influence behavior will need specialist support and tooling. This might require tools embedded within the user interface to evaluate and monitor user interaction, and capture user feedback(56), or directed implementation studies.(57)

Despite the oft cited risks of alert fatigue with Clinical Decision Support Systems (CDSS)(58), there is good evidence that well designed alerts can be impactful.(59, 60, 53) Overt behavioural modifications will need a mechanism to explain their recommendation (as per XAI) or generate trust (see Pillar 5).(61) Trust will possibly be more important where behavior modification is indirect through non-interruptive techniques (e.g. re-ordering preference lists or otherwise adapting the user interface to make the recommended choice more accessible).

5 CONTINUOUS CLINICAL EVALUATION

Our analogy with translational medicine breaks down at the evaluation stage. For drug discovery, evaluation is via a randomized controlled trial (RCT). Randomization handles unanticipated bias and ML4H should hold itself to the same standard but of 350,000 studies registered on ClinicalTrials.gov⁷ in 2020, just 358 evaluated ML4H, and only 66 were randomized.(62) As usual for ML4H, those RCTs were not interacting with EHR data. They were evaluations of algorithms supporting imaging, cataract screening, colonoscopy, cardiocographs and more.(63, 64, 65, 66, 67, 68, 69)

Where the ML4H intervention delivers a novel biological treatment strategy, then it is appropriate to reach for the full paraphernalia used in Clinical Trials of Investigational Medicinal Products (CTIMPs).(2) But in many cases, algorithms will be used to optimize operational workflows and clinical pathways. These pathways may be specific and contextual rather than generalizable. Poor external validity is not a critique: an algorithm that is useful or important in one institution does not have to be relevant in the next (the 'myth of generalizability').(70) Moreover, the algorithm is not the same as the patented and fixed active ingredient in a medicinal product. This is no single point in time nor single host environment at which it can be declared enduringly effective. This means that institutions deploying and relying on these tools need a strategy for rapid continuous clinical and operational evaluation.

This time the EHR may provide an advantage instead of just additional complexity. Since ML4H algorithms must be implemented through some form of direct or indirect CDSS, then the next logical step is to randomize the deployment of those alerts. This in itself is not novel. Randomized deterministic alerts from CDSS are part of the standard evaluation toolkit for quality improvement initiatives in at NYU

⁷ <ClinicalTrials.gov>

Langone(71), and for research elsewhere.(72) At NYU Langone, such tooling permitted a small team to deliver 10 randomized trials within a single year.(71)

The final pillar in our CDE uses the same approach for the probabilistic insights derived from ML4H. Excellent patient and public involvement, and ethical guidance, will be required to distinguish those algorithms that require per patient point-of-care consent from those that can use opt-out or cluster methods. But we think that latter group is large for two reasons. Firstly, patients are exposed to varying treatment regimes by dint of their random interaction with different clinicians based on geography (the healthcare provider they access) and time (staff holidays and shift patterns etc.). This routine variation in practice is summarized as the 60-30-10 problem: 60% of care follows best practice; 30% is wasteful or ineffective and 10% is harmful.(6) Secondly, because the intervention is informational, there is ethical precedent for patient level randomization without consent (e.g. Acute Kidney Injury alerts).(72) This hints at a larger and more routine role for randomization in evaluation of algorithms. This in turn is supported by a growing(52, 73, 74) but sometimes conflicting(75) literature on opt-out consent in Learning Healthcare Systems (LHS). As such, progress will require careful attention to a range of concerns.

At our own institution, we have extended this ethical and safety case one step further, and we are piloting a study design where the randomization is non-mandatory: a nudge not an order.(76) The clinician is explicitly invited to only comply with the randomization where they have equipoise themselves. Where they have a preference, they overrule the alert (see Vignette 1 in the Electronic Supplemental Material).

Embedded randomized digital evaluation should permit rapid evidence generation, and build the trust needed to support the implementation described under Pillar 4.

DRUG DISCOVERY PARALLELS

We have described a template for a Clinical Deployment Environment that supports the translation of ML4H algorithms from bench to bedside. Although the requirements differ, the objective is similar to that for drug development. A similar approach to phasing has previously been proposed for (biomarker) prediction models.(77)

Most ML4H that derives value from the EHR is in the pre-clinical phase. In drug development, the objective of this phase is to identify candidate molecules which might make effective drugs. Evaluation is conducted in vitro. Metrics used to evaluate candidates, such binding affinity or other pharmacokinetic properties, describe the properties of the molecule.(78) For ML, the objective is to identify candidate algorithms, comprising of input variables and model structures, which might make the core of an effective CDSS. Evaluation is conducted offline on de-identified datasets. Metrics used to evaluate candidates, such Area Under the Receiver Operator Curve (AUROC), the F1 score and calibration, describe the properties of the algorithm.(79)

Phase 1 drug trials are the first time a drug candidate is tested in humans. They are conducted in small numbers of healthy volunteers. The aim of the trial is to determine the feasibility of progressing to trials in patients by determining drug safety and appropriate dosage. Drug formulation, the processes by which substances are combined with the active pharmaceutical ingredient to optimize the acceptability and effective delivery of the drug, is also considered at this stage. Phase 1 ML4H trials are the first time an algorithm candidate is tested within the healthcare environment. The aim of the trial is to determine the feasibility of progressing to trials of efficacy by ensuring the algorithm implementation is safe, reliable and able to cope with real-world data quality issues. The development of a mechanism to deliver of algorithm outputs embedded in the clinical workflow is also be considered at this stage.

Phase 2 drug trials involve recruitment of small numbers patients with the disease of interest, typically 50 – 200. The aim is to determine drug efficacy at treating the disease. Treating clinicians are involved in so far as they must agree to prescribe the drug for their patients. The trials are often too short to determine long term outcomes, therefore surrogate measures such biomarker status or change in tumour size are used as endpoints.(80) Phase 2 ML4H trials involve recruitment of small numbers of clinicians making the decision of interest, typically 5 – 10. The aim is to determine the efficacy of the algorithm in improving their decisions. Patients are involved in so far as they must agree to be on the receiving end of these supported decisions and identifiable data is required. Endpoints are markers of successful task completion in all cases. Investigations to determine ways in which the system could be more successful in influencing user behavior are carried out at this stage. These include usability analyses, considerations of how well the ML4H/CDSS is integrated into the overall system and implementation studies to identify how best to optimize end-user adoption and engagement.(57)

Phase 3 drug trials involve the recruitment of large numbers of patients to determine whether a drug is effective in improving patient outcomes. The gold standard of trial design is a double-blinded randomized controlled trial (RCT). Phase 3 ML4H trials will require integration of data from multiple centers for algorithms acting on specific decisions but inevitably adapted to their local data environment.

The phases of drug development are not meant to be matched 1:1 to the pillars of the CDE described here: in fact, our argument for 'real-world' *development* deliberately seeks to merge the steps. But the parallel is drawn to highlight the effort necessary to see ML4H have an impact on the clinical and operational decision making in the workplace. Heretofore this effort has been hugely underestimated.

CONCLUSION

Even this analogy stops short of the full task of deployment. With drug development, the universities and the pharmaceutical industry go on to take advantage of a supply chain to deliver the drug to the hospital with the necessary quality control and monitoring. Those prescribing and administering the drug have spent years in training, and are supported by pharmacists and medication safety experts. And even after the drug is administered, observation and long term follow-up continue to identify side-effects and long term hazards.

That network of expertise and infrastructure is largely in place where software *is within* (not *as*) a medical device, but is only just being envisioned where the data driving ML4H comes from the EHR. This distinction needs to be made else the disillusionment with the promise of ML4H will continue. The technology does have the potential to change how we deliver health but the methodology alone is insufficient. The impressive demonstrations of the power of AI and ML to beat humans in games, and predict protein structures does not mean that these tools are ready for wide spread deployment.

But we should not be pessimistic. As per author William Gibson, it is clear that "The Future Has Arrived — It's Just Not Evenly Distributed." Beyond healthcare, machine learning has already demonstrated that it can reliably create value.(5) It is now our responsibility to take those lessons and adapt them for our patients.

The Five Pillars outlined here are a sketch of that redistribution. They are born from our local experience (Pillars 1, 2 and 5) and our wider observations (Pillars 3 and 4). They fundamentally are an argument for a professionalization of ML4H, and a caution against the 'get-rich quick' headlines in the popular and scientific press.(1) We envision a future where each algorithm is managed in a digital pharmacy with the same rigor that we apply to medicines. But unlike drugs, some of these algorithms will have their entire

life-cycle, from development to deployment, managed by the local healthcare provider. Computer vision tasks that support diagnostic radiology can be partially developed offline. Components of sepsis prediction tools will transfer from institution to institution but will need adapting to local clinical workflows. But there will be opportunity and value for ML4H to optimize operational tasks that are temporary or specific to that institution. This means that some development and much of the deployment will require a suitably trained workforce, and an infrastructure perhaps supported by these five pillars.

CONFLICT OF INTEREST STATEMENT

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

AUTHOR CONTRIBUTIONS

All authors listed have made a substantial, direct and intellectual contribution to the work, and approved it for publication.

FUNDING

SH is supported by a Health Foundation Improvement Science Fellowship, and by funding from the Biomedical Research Centre based at the University College Hospitals (UCLH) National Health Service (NHS) Foundation Trust and University College London (UCL). NS is supported by funding from Health Data Research UK. SH, TK, and NS are supported funds from the National Institute for Health Research (Artificial Intelligence, Digitally adapted, hyper-local real-time bed forecasting to manage flow for NHS wards, AI AWARD01786) and NHS-X. WL is supported by a UKRI Ernest Rutherford Fellowship. The views expressed in this publication are those of the author(s) and not necessarily those of the National Institute for Health Research, NHS-X or the Department of Health and Social Care.

REFERENCES

- 1 Bunz M, Braghieri M. The AI doctor will see you now: assessing the framing of AI in news coverage. *AI & SOCIETY* **37** (2022) 9–22. doi:10.1007/s00146-021-01145-9.
- 2 Komorowski M, Celi LA, Badawi O, Gordon AC, Faisal AA. The Artificial Intelligence Clinician learns optimal treatment strategies for sepsis in intensive care. *Nature Medicine* **24** (2018) 1716–1720. doi:10.1038/s41591-018-0213-5.
- 3 Hyland SL, Faltys M, Hüser M, Lyu X, Gumbsch T, Esteban C, et al. Early prediction of circulatory failure in the intensive care unit using machine learning. *Nature Medicine* **26** (2020) 364–373.
- 4 Tomašev N, Glorot X, Rae J, Zielinski M, Askham H, Saraiva A, et al. A clinically applicable approach to continuous prediction of future acute kidney injury. *Nature* **572** (2019) 116–119.
- 5 McRae H. Facebook, Airbnb, Uber, and the unstoppable rise of the content non-generators. *The Independent* (2015).
- 6 Braithwaite J, Glasziou P, Westbrook J. The three numbers you need to know about healthcare: the 60-30-10 challenge. *BMC medicine* **18** (2020) 1–8.
- 7 The national strategy for AI in health and social care (2022).
- 8 Digital future index 2021-2022 (2021).

- 322 **9** Everson J, Rubin JC, Friedman CP. Reconsidering hospital EHR adoption at the dawn of HITECH:
 323 implications of the reported 9% adoption of a “basic” EHR. *Journal of the American Medical Informatics*
 324 *Association* **27** (2020) 1198–1205. doi:10.1093/jamia/ocaa090.
- 325 **10** Muehlematter U, Daniore P, Vokinger K. Approval of artificial intelligence and machine learning-based
 326 medical devices in the USA and Europe (2015–20): a comparative analysis. *The Lancet. Digital health*
 327 **3** (2021) e195–e203.
- 328 **11** Morse K, Bagley S, Shah N. Estimate the hidden deployment cost of predictive models to improve
 329 patient care. *Nature medicine* **26** (2020) 18–19.
- 330 **12** Eaneff S, Obermeyer Z, Butte AJ. The Case for Algorithmic Stewardship for Artificial Intelligence and
 331 Machine Learning Technologies. *JAMA* **324** (2020) 1397–1398. doi:10.1001/jama.2020.9371.
- 332 **13** Burton PR, Murtagh MJ, Boyd A, Williams JB, Dove ES, Wallace SE, et al. Data Safe Havens in
 333 health research and healthcare. *Bioinformatics (Oxford, England)* **31** (2015) 3241–3248. doi:10.1093/
 334 bioinformatics/btv279.
- 335 **14** Keane PA, Topol EJ. With an eye to AI and autonomous diagnosis. *npj Digital Medicine* **1** (2018) 40,
 336 s41746–018–0048–y. doi:10.1038/s41746-018-0048-y.
- 337 **15** Woolf SH. The Meaning of Translational Research and Why It Matters. *JAMA* **299** (2008). doi:10.
 338 1001/jama.2007.26.
- 339 **16** Corrigan-Curay J, Sacks L, Woodcock J. Real-World Evidence and Real-World Data for Evaluating
 340 Drug Safety and Effectiveness. *JAMA* **320** (2018) 867. doi:10.1001/jama.2018.10136.
- 341 **17** Williamson E, Walker A, Bhaskaran K, Bacon S, Bates C, Morton C, et al. Factors associated with
 342 COVID-19-related death using OpenSAFELY. *Nature* **584** (2020) 430–436.
- 343 **18** Data research infrastructure landscape: A review of the UK data research infrastructure (2021).
- 344 **19** The DECIDE-AI Steering Group. DECIDE-AI: new reporting guidelines to bridge the development-to-
 345 implementation gap in clinical artificial intelligence. *Nature Medicine* **27** (2021) 186–187. doi:10.1038/
 346 s41591-021-01229-5.
- 347 **20** Shah NH, Milstein A, Bagley SC PhD. Making Machine Learning Models Clinically Useful. *JAMA*
 348 **322** (2019) 1351. doi:10.1001/jama.2019.10306.
- 349 **21** Guinney J, Saez-Rodriguez J. Alternative models for sharing confidential biomedical data. *Nature*
 350 *biotechnology* **36** (2018) 391–392.
- 351 **22** DevOps (2022).
- 352 **23** Ben-Israel D, Jacobs WB, Casha S, Lang S, Ryu WHA, de Lotbiniere-Bassett M, et al. The impact of
 353 machine learning on patient care: A systematic review. *Artificial Intelligence in Medicine* **103** (2020)
 354 101785.
- 355 **24** Sculley D, Holt G, Golovin D, Davydov E, Phillips T, Ebner D, et al. Hidden technical debt in machine
 356 learning systems. *Advances in neural information processing systems* **28** (2015) 2503–2511.
- 357 **25** SEC Charges Knight Capital With Violations of Market Access Rule (2013).
- 358 **26** John MM, Olsson HH, Bosch J. Towards MLOps: A framework and maturity model. *Towards MLOps: A*
 359 *framework and maturity model* (IEEE), vol. 2021 47th Euromicro Conference on Software Engineering
 360 and Advanced Applications (SEAA) (2021), 1–8.
- 361 **27** Feng J, Phillips RV, Malenica I, Bishara A, Hubbard AE, Celi LA, et al. Clinical artificial intelligence
 362 quality improvement: towards continual monitoring and updating of AI algorithms in healthcare. *npj*
 363 *Digital Medicine* **5** (2022) 66. doi:10.1038/s41746-022-00611-y.
- 364 **28** Artificial Intelligence and Machine Learning (AI/ML) Software as a Medical Device Action Plan. Tech.
 365 rep., US Food & Drug Administration (2021).

- 366 **29** Feng J, Emerson S, Simon N. Approval policies for modifications to machine learning-based software
367 as a medical device: A study of bio-creep. *Biometrics* **77** (2021) 31–44. doi:10.1111/biom.13379.
- 368 **30** Amodei D, Olah C, Steinhardt J, Christiano P, Schulman J, Mané D. Concrete problems in AI safety.
369 *arXiv preprint arXiv:1606.06565* (2016). doi:10.48550/arXiv.1606.06565.
- 370 **31** Liley J, Emerson S, Mateen B, Vallejos C, Aslett L, Vollmer S. Model updating after interventions
371 paradoxically introduces bias. *Proceedings of Machine Learning Research* (2021), 3916–3924.
- 372 **32** Gunning D, Stefik M, Choi J, Stumpf S, Yang G. XAI—Explainable artificial intelligence. *Science*
373 *Robotics* **4** (2019).
- 374 **33** Mueller S, Hoffman R, Clancey W, Emrey A, Klein G. Explanation in human-AI systems: A literature
375 meta-review, synopsis of key ideas and publications, and bibliography for explainable AI. *arXiv preprint*
376 *arXiv:1902.01876* (2019).
- 377 **34** Vilone G, Longo L. Explainable artificial intelligence: a systematic review. *arXiv preprint*
378 *arXiv:2006.00093* (2020).
- 379 **35** Linardatos P, Papastefanopoulos V, Kotsiantis S. Explainable ai: A review of machine learning
380 interpretability methods. *Entropy. An International and Interdisciplinary Journal of Entropy and*
381 *Information Studies* **23** (2020).
- 382 **36** Ghassemi M, Oakden-Rayner L, Beam A. The false hope of current approaches to explainable artificial
383 intelligence in health care. *The Lancet. Digital health* **3** (2021) e745–e750.
- 384 **37** The Disagreement Problem in Explainable Machine Learning: A Practitioner’s Perspective (2022).
385 ArXiv:2202.01602 [cs].
- 386 **38** Doshi-Velez F, Kim B. Towards a rigorous science of interpretable machine learning. *arXiv preprint*
387 *arXiv:1702.08608* (2017).
- 388 **39** Hoffman R, Mueller S, Klein G, Litman J. Metrics for explainable AI: Challenges and prospects. *arXiv*
389 *preprint arXiv:1812.04608* (2018).
- 390 **40** Abdar M, Pourpanah F, Hussain S, Rezazadegan D, Liu L, Ghavamzadeh M, et al. A review of
391 uncertainty quantification in deep learning: Techniques, applications and challenges. *Information*
392 *Fusion* **76** (2021) 243–297.
- 393 **41** Leibig C, Allken V, Ayhan M, Berens P, Wahl S. Leveraging uncertainty information from deep neural
394 networks for disease detection. *Scientific reports* **7** (2017) 1–14.
- 395 **42** Filos A, Farquhar S, Gomez A, Rudner T, Kenton Z, Smith L, et al. A systematic comparison of
396 bayesian deep learning robustness in diabetic retinopathy tasks. *arXiv preprint arXiv:1912.10481*
397 (2019).
- 398 **43** Ghoshal B, Tucker T. Estimating uncertainty and interpretability in deep learning for coronavirus
399 (COVID-19) detection. *arXiv preprint arXiv:2003.10769* (2020).
- 400 **44** Chow C. On optimum recognition error and reject tradeoff. *IEEE Transactions on Information Theory*
401 **16** (1970) 41–46. doi:10.1109/TIT.1970.1054406.
- 402 **45** Bartlett PL, Wegkamp MH. Classification with a Reject Option using a Hinge Loss. *Journal of Machine*
403 *Learning Research* **9** (2008) 18.
- 404 **46** Tortorella F. An Optimal Reject Rule for Binary Classifiers. Goos G, Hartmanis J, van Leeuwen J, Ferri
405 FJ, Iñesta JM, Amin A, et al., editors, *Advances in Pattern Recognition* (Berlin, Heidelberg: Springer
406 Berlin Heidelberg), vol. 1876 (2000), 611–620. doi:10.1007/3-540-44522-6_63. Series Title: Lecture
407 Notes in Computer Science.
- 408 **47** El-Yaniv R, Wiener Y. On the Foundations of Noise-free Selective Classification. *Journal of Machine*
409 *Learning Research* **11** (2010) 37.

- 410 **48** Davis SE, Lasko TA, Chen G, Siew ED, Matheny ME. Calibration drift in regression and machine
411 learning models for acute kidney injury. *Journal of the American Medical Informatics Association* **24**
412 (2017) 1052–1061. doi:10.1093/jamia/ocx030.
- 413 **49** Davis SE, Greevy RA, Lasko TA, Walsh CG, Matheny ME. Detection of calibration drift in clinical
414 prediction models to inform model updating. *Journal of Biomedical Informatics* **112** (2020) 103611.
415 doi:10.1016/j.jbi.2020.103611.
- 416 **50** Parisi G, Kemker R, Part J, Kanan C, Wermter S. Continual lifelong learning with neural networks: A
417 review. *Neural Networks* **113** (2019) 54–71.
- 418 **51** Lee CS, Lee AY. Clinical applications of continual learning machine learning. *The Lancet Digital*
419 *Health* **2** (2020) e279–e281. doi:10.1016/S2589-7500(20)30102-3. Publisher: Elsevier.
- 420 **52** Halpern SD. Using Default Options and Other Nudges to Improve Critical Care:. *Critical Care Medicine*
421 **46** (2018) 460–464. doi:10.1097/CCM.0000000000002898.
- 422 **53** Main C, Moxham T, Wyatt JC, Kay J, Anderson R, Stein K. Computerised decision support systems
423 in order communication for diagnostic, screening or monitoring test ordering: systematic reviews of
424 the effects and cost-effectiveness of systems. *Health Technology Assessment* **14** (2010). doi:10.3310/
425 hta14480.
- 426 **54** Sendak M, Ratliff W, Sarro D, Alderton E, Futoma J, Gao M, et al. Real-world integration of a sepsis
427 deep learning technology into routine clinical care: Implementation study. *JMIR medical informatics* **8**
428 (2020) e15182.
- 429 **55** Connell A, Black G, Montgomery H, Martin P, Nightingale C, King D, et al. Implementation of
430 a Digitally Enabled Care Pathway (Part 2): Qualitative Analysis of Experiences of Health Care
431 Professionals. *Journal of Medical Internet Research* **21** (2019) e13143. doi:10.2196/13143.
- 432 **56** Yusop NSM, Grundy J, Vasa R. Reporting Usability Defects: A Systematic Literature Review. *IEEE*
433 *Transactions on Software Engineering* **43** (2017) 848–867. doi:10.1109/TSE.2016.2638427. Conference
434 Name: IEEE Transactions on Software Engineering.
- 435 **57** Sutton RT, Pincock D, Baumgart DC, Sadowski DC, Fedorak RN, Kroeker KI. An overview of clinical
436 decision support systems: benefits, risks, and strategies for success. *npj Digital Medicine* **3** (2020) 1–10.
437 doi:10.1038/s41746-020-0221-y. Number: 1 Publisher: Nature Publishing Group.
- 438 **58** Phansalkar S, van der Sijs H, Tucker A, Desai A, Bell D, Teich J, et al. Drug-drug interactions that
439 should be non-interruptive in order to reduce alert fatigue in electronic health records. *Journal of the*
440 *American Medical Informatics Association : JAMIA* **20** (2013) 489–493.
- 441 **59** Park G, Kang B, Kim S, Lee J. Retrospective review of missed cancer detection and its mammography
442 findings with artificial-intelligence-based, computer-aided diagnosis. *Diagnostics (Basel, Switzerland)*
443 **12** (2022) 387.
- 444 **60** Sayres R, Taly A, Rahimy E, Blumer K, Coz D, Hammel N, et al. Using a Deep Learning Algorithm
445 and Integrated Gradients Explanation to Assist Grading for Diabetic Retinopathy. *Ophthalmology* **126**
446 (2019) 552–564. doi:10.1016/j.optha.2018.11.016.
- 447 **61** McCoy L, Brenna C, Chen S, Vold K, Das S. Believing in black boxes: machine learning for healthcare
448 does not need explainability to be evidence-based. *Journal of clinical epidemiology* **142** (2022) 252–257.
- 449 **62** Zippel C, Bohnet-Joschko S. Rise of clinical studies in the field of machine learning: A review of data
450 registered in ClinicalTrials.gov. *International journal of environmental research and public health* **18**
451 (2021) 5072.
- 452 **63** INFANT CG. Computerised interpretation of fetal heart rate during labour (INFANT): a randomised
453 controlled trial. *Lancet (London, England)* **389** (2017) 1719–1729.

- 454 **64** Titano JJ, Badgeley M, Schefflein J, Pain M, Su A, Cai M, et al. Automated deep-neural-network
455 surveillance of cranial images for acute neurologic events. *Nature Medicine* **24** (2018) 1337–1341.
456 doi:10.1038/s41591-018-0147-y.
- 457 **65** Wang P, Berzin T, Glissen Brown J, Bharadwaj S, Becq A, Xiao X, et al. Real-time automatic detection
458 system increases colonoscopic polyp and adenoma detection rates: a prospective randomised controlled
459 study. *Gut* **68** (2019) 1813–1819.
- 460 **66** Wu L, Zhang J, Zhou W, An P, Shen L, Liu J, et al. Randomised controlled trial of WISENSE, a
461 real-time quality improving system for monitoring blind spots during esophagogastroduodenoscopy.
462 *Gut* **68** (2019) 2161–2169.
- 463 **67** Lin H, Li R, Liu Z, Chen J, Yang Y, Chen H, et al. Diagnostic efficacy and therapeutic decision-making
464 capacity of an artificial intelligence platform for childhood cataracts in eye clinics: A multicentre
465 randomized controlled trial. *EClinicalMedicine* **9** (2019) 52–59.
- 466 **68** Turakhia MP, Desai M, Hedlin H, Rajmane A, Talati N, Ferris T, et al. Rationale and design of a
467 large-scale, app-based study to identify cardiac arrhythmias using a smartwatch: The Apple Heart Study.
468 *American Heart Journal* **207** (2019) 66–75. doi:10.1016/j.ahj.2018.09.002.
- 469 **69** Long E, Lin H, Liu Z, Wu X, Wang L, Jiang J, et al. An artificial intelligence platform for the
470 multihospital collaborative management of congenital cataracts. *Nature Biomedical Engineering* **1**
471 (2017).
- 472 **70** Futoma J, Simons M, Panch T, Doshi-Velez F, Celi LA. The myth of generalisability in clinical research
473 and machine learning in health care. *The Lancet Digital Health* **2** (2020) e489–e492. doi:10.1016/
474 S2589-7500(20)30186-2.
- 475 **71** Horwitz LI, Kuznetsova M, Jones SA. Creating a Learning Health System through Rapid-Cycle,
476 Randomized Testing. *New England Journal of Medicine* **381** (2019) 1175–1179. doi:10.1056/
477 NEJMs1900856.
- 478 **72** Wilson FP, Martin M, Yamamoto Y, Partridge C, Moreira E, Arora T, et al. Electronic health record
479 alerts for acute kidney injury: multicenter, randomized clinical trial. *BMJ* (2021) m4786. doi:10.1136/
480 bmj.m4786.
- 481 **73** London AJ. Learning health systems, clinical equipoise and the ethics of response adaptive
482 randomisation. *Journal of Medical Ethics* **44** (2018) 409–415. doi:10.1136/medethics-2017-104549.
- 483 **74** Scobie S, Castle-Clarke S. Implementing learning health systems in the UK NHS: Policy actions to
484 improve collaboration and transparency and support innovation and better use of analytics. *Learning*
485 *Health Systems* **4** (2020) e10209. doi:10.1002/lrh2.10209.
- 486 **75** Meyer MN, Heck PR, Holtzman GS, Anderson SM, Cai W, Watts DJ, et al. Objecting to experiments
487 that compare two unobjectionable policies or treatments. *Proceedings of the National Academy of*
488 *Sciences* **116** (2019) 10723–10728. doi:10.1073/pnas.1820701116.
- 489 **76** Wilson MG, Asselbergs FW, Harris SK. Learning from individualised variation for evidence generation
490 within a learning health system. *British Journal of Anaesthesia* (2022).
- 491 **77** Pepe MS, Etzioni R, Feng Z, Potter JD, Thompson ML, Thornquist M, et al. Phases of Biomarker
492 Development for Early Detection of Cancer. *JNCI: Journal of the National Cancer Institute* **93** (2001)
493 1054–1061. doi:10.1093/jnci/93.14.1054.
- 494 **78** Lo YC, Rensi SE, Torng W, Altman RB. Machine learning in chemoinformatics and drug discovery.
495 *Drug Discovery Today* **23** (2018) 1538–1546. doi:10.1016/j.drudis.2018.05.010.
- 496 **79** Collins GS, Reitsma JB, Altman DG, Moons KGM. Transparent Reporting of a multivariable prediction
497 model for Individual Prognosis Or Diagnosis (TRIPOD): The TRIPOD Statement. *Annals of Internal*
498 *Medicine* **162** (2015) 55–11. doi:10.7326/M14-0697.

499 **80** Van Norman GA. Phase II Trials in Drug Development and Adaptive Trial Design. *JACC: Basic to*
 500 *Translational Science* 4 (2019) 428–437. doi:10.1016/j.jacbs.2019.02.005.

FIGURES

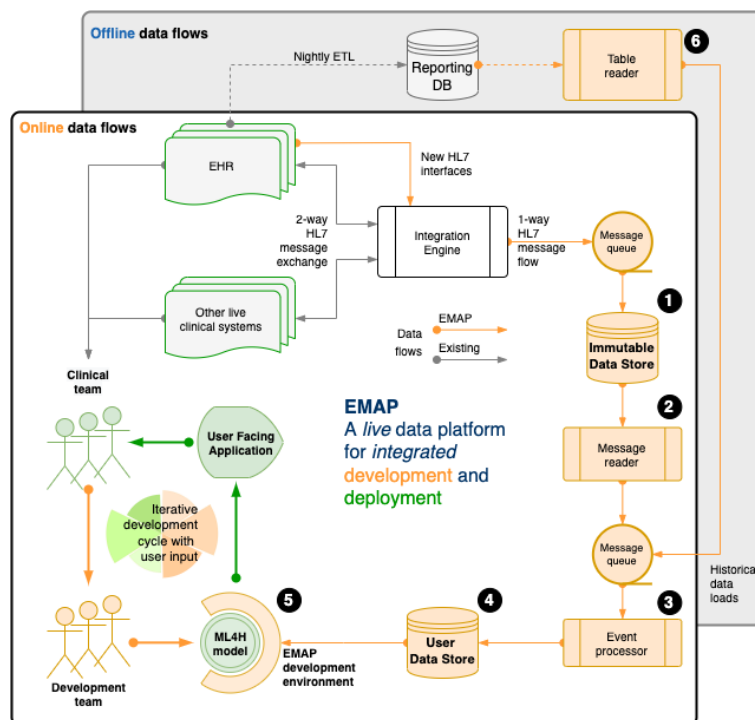


Figure 1. Our real-world development is performed on the Experimental Medicine Application Platform (EMAP). EMAP is a clinical laboratory within which ML4H researchers can iteratively build, test and gather feedback from the bedside. It unifies the data and the tools for off-line and online development of ML4H models (see Figure 1 and the **(numbers)** in the following sentences that refer to objects in the figure).

In brief, EMAP builds a patient orientated SQL database from Health Level 7 version 2 (HL7v2) messages that are being exchanged between hospital systems. HL7v2 messages are ubiquitous in health care, and the *de facto* standard for internal communication. Rather than multiple pairwise connections between different hospital electronic systems, an integration engine acts as a single hub that routes HL7 messages, and where necessary translates to ensure compatibility. EMAP copies each message passing through the integration engine to a PostgreSQL database, the *Immutable Data Store* (IDS) (1). A *message reader* (2) processes each live message to an interchange format so that downstream processing is insulated from local HL7 implementation. Separately, the *table reader* (6) processes historical data (e.g. from the reporting database) to the same interchange format. Live messages take priority over historical messages in a queue that feeds the *event processor* (3). This links each message to a patient and a hospital visit, makes appropriate updates for out of order messages, and merges when separate identifiers are recognised to represent the same patient. A full audit trail is maintained. Each event updates a second live PostgreSQL database, the *User Data Store* (UDS) (4). The hospital hosts Jupyter and RStudio servers, and a Linux development environment is provided that allows docker deployment, installation of analysis libraries and frameworks, exposes SSH and HTTPS services, and allows user verification against the hospital active directory. (5) A typical workflow might include investigation and experimentation in a Jupyter Notebook with data from the UDS, then using a small network of docker containers to run the development script, log outputs to a testing database, and report to users via email or a locally hosted web application or dashboard. A fuller explanation is available in the Electronic Supplementary Material (Section 2: EMAP data flows).

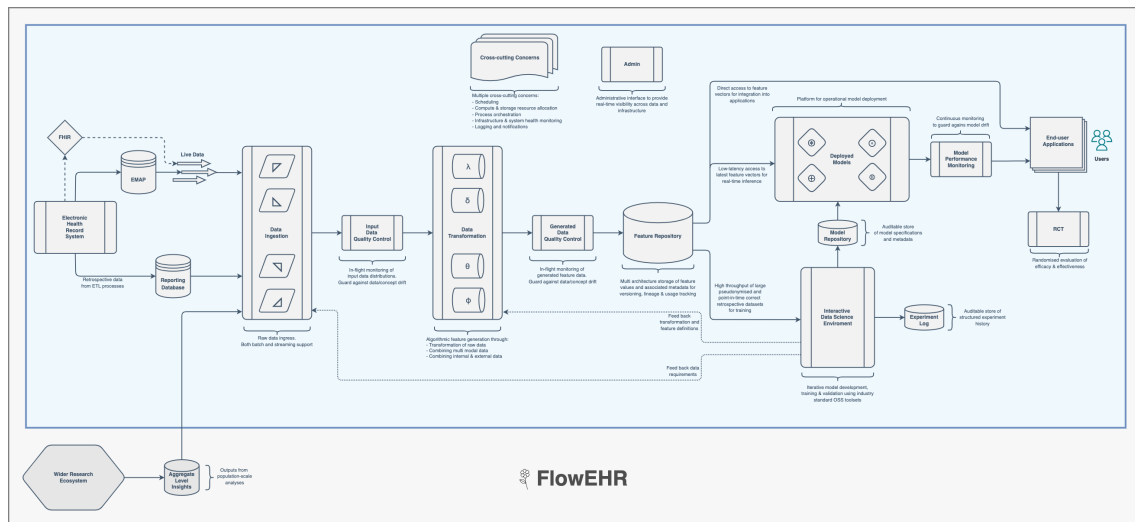


Figure 2. Our ML-Ops platform is called FlowEHR. Moving from left to right across the figure, the system monitors raw input data including checks for distribution shift, builds features with testable and quality controlled code, makes those features available for both training and predictions to avoid train/serve skew, and maintains an auditable and monitored model repository.