# MORF: A Framework for Predictive Modeling and Replication At Scale With Privacy-Restricted MOOC Data

Josh Gardner, Christopher Brooks
*School of Information*
*The University of Michigan*
*Ann Arbor, USA*
*{jpgard, brooksch}@umich.edu*

Juan Miguel Andres, Ryan S. Baker
*Graduate School of Education*
*The University of Pennsylvania*
*Philadelphia, USA*
*miglimjapandres@gmail.com, rybaker@upenn.edu*

*Abstract*—Big data repositories from online learning platforms such as Massive Open Online Courses (MOOCs) represent an unprecedented opportunity to advance research on education at scale and impact a global population of learners. To date, such research has been hindered by poor reproducibility and a lack of replication, largely due to three types of barriers: experimental, inferential, and data. We present a novel system for large-scale computational research, the MOOC Replication Framework (MORF), to jointly address these barriers. We discuss MORF's architecture, an open-source platform-as-a-service (PaaS) which includes a simple, flexible software API providing for multiple modes of research (predictive modeling or production rule analysis) integrated with a high-performance computing environment. All experiments conducted on MORF use executable Docker containers which ensure complete reproducibility while allowing for the use of any software or language which can be installed in the linux-based Docker container. Each experimental artifact is assigned a DOI and made publicly available. MORF has the potential to accelerate and democratize research on its massive data repository, which currently includes over 200 MOOCs, as demonstrated by initial research conducted on the platform. We also highlight ways in which MORF represents a solution template to a more general class of problems faced by computational researchers in other domains.

*Keywords*-data infrastructure; education; MOOC; reproducibility; predictive modeling; machine learning

## I. INTRODUCTION

### A. Educational Big Data in the MOOC Era

Since the launch of the three major Massive Open Online Course (MOOC) platforms in 2012, MOOCs have been met variously with enthusiasm, curiosity, and criticism regarding their potential to transform the educational landscape and provide unprecedented access to high-quality educational content and credentials for learners around the globe. The massive, granular, and multimodal data generated by MOOCs has the potential to answer many of the research questions underlying these perspectives, informing educational research on the process of learning within online environments. To date, MOOC research has addressed a diverse array of research questions from psychometrics and social psychology to predictive modeling and machine learn-

ing. For example, several works have explored prediction of various student outcomes using behavioral, linguistic, and assignment data from MOOCs to evaluate and predict various student outcomes including course completion [1], [2], [3], assignment grades [4], Correct on First Attempt (CFA) submissions [5], student confusion [6], and changes in behavior over time [7]. A key area of research has been methods for feature engineering, or extracting structured information from raw data (i.e. clickstream server logs, natural language in discussion posts) [8].

### B. A Confluence of Challenges for MOOC Research

The challenges of doing big data research have grown over the past decade as the data, statistical models, and technical tools have become increasingly complex, and MOOC research has not been immune to the many challenges facing the larger big data research community. Issues with reproducibility and the lack of replication studies have the potential to hamper the field and limit researchers' understanding, and, consequently, to diminish the potential impact of MOOCs for learners around the globe. Strict privacy regulations also limit the sharing of learner data for MOOC research. Many universities interpret the Family Educational Rights and Privacy Act (FERPA), the IRB Common Rule, or other data regulations in ways that severely limit access to MOOC data. The recent passage of the General Data Protection Regulation (GDPR) in the European Union may further limit data sharing and collection.

### C. Contribution: A Proposed Solution

In this work, we present a novel big data research platform, the MOOC Replication Framework (MORF), which is designed to jointly address technical, data, and methodological barriers to reproducible and replication research in MOOCs. We present empirical evidence on the current state of the "replication crisis" in big data and the learning sciences in Section II, motivating the need for MORF. Then, we provide a detailed overview of the MORF platform in Section III, detailing MORF's unique platform-as-a-service (PaaS) architecture for working with massive

privacy-restricted MOOC datasets in two modes of analysis (predictive modeling and production rule analysis). In particular, we highlight its use of containerization to ensure full replicability of computational results, and its massive, highly diverse datasets available for analysis in a secure, high-performance environment. We discuss how MORF represents a framework which can apply broadly to many other research domains, the software's openness and reuse potential, and the many benefits of big data replication research, in Section IV.

## II. THE REPLICATION CRISIS: AN EMPIRICAL PERSPECTIVE

An emerging body of empirical evidence is revealing the extent to which several factors – a lack of basic reproducibility practices in big data research; a dearth of replication studies; and experimental, data, and inferential challenges to big data research – have hindered the field.

### A. Reproducibility in Big Data Modeling Research

Recent evidence has demonstrated that issues with basic reproducibility – the ability to regenerate the original results of an experiment using the original methods and data – are widespread in the field of big data research, which covers a variety of methods spanning artificial intelligence, machine learning, data mining, and simulation-based research. In a survey of 400 papers from leading artificial intelligence venues, none documented all aspects necessary to fully reproduce the work; on average, only 20-30% of the factors needed to replicate the original work were reported [9]. In a survey of 30 text mining studies, none of the 30 works surveyed provided source code, and only one of 16 applicable studies provided an executable, to reproduce their experiments [10]; lack of access to data, software environment, and implementation methods were noted as barriers to reproducibility in the works surveyed. Even when code is available, it is often insufficient to reproduce an experiment: in a survey of 613 published computer systems papers, the published code accompanying failed to build or run in 20% of cases; in total, it was not possible to verify or reproduce 75.1% of studies surveyed using the artifacts provided in publication [11].

Much-needed replication studies – where the original methods of an experiment are applied to *new* data to evaluate the generalizability of findings and contextualize the results – are not even possible with a great deal of published big data research, as even mere reproducibility is not possible.

### B. Lack of Replication Studies in Learning Sciences

Replication research in the domain of the learning sciences has been scarce to date. A survey of the top 100 education journals [12] demonstrated that only 0.13% of education articles were replications (a replication rate eight times lower than in the field of psychology). 67.4% of replications attempted replicated the original results fully, 19.5% replicated some (but not all) findings, and 13.1% failed to replicate any original findings.

Replication specific to MOOCs and educational big data research has been particularly scarce. This may make the existing body of research in MOOCs particularly unreliable, an issue compounded by the fact that most MOOC studies to date have used small samples of MOOCs, and selected highly-varying subsets of students from the available datasets (e.g. only students who joined in the first 10 days of the course and have viewed at least one video; completed pre-course survey and the first end-of-unit exam, etc.) [8]. In other domains, it has been shown that selecting subpopulations which vary in even subtle ways has led to very different experimental conclusions in big data analyses depending on the population used [13].

The limited MOOC replication research to date has shown that many published findings in the field may not replicate. For example, in a large-scale replication, [14] found that only 12 of 15 previously-published MOOC findings replicated significantly across the data sets, and that two findings replicated significantly in the opposite direction; similar results in a machine learning replication in the context of algorithm and feature selection were shown in [15].

### C. Three Barriers to Reproducible Big Data Research

We identify three key sets of barriers contributing to both the lack of general reproducibility in big data research, and the lack of replication studies within the field of educational big data research in particular. MORF attempts to resolve all three barriers.

**Experimental challenges** with reproducibility relate to reproducing the exact experimental protocol [9]. Many have advocated for the open sharing of code as a potential solution to address technical issues with reproducibility [16]. However, as discussed in Section II-A, even when code is available, other technical issues can prevent reproducibility in computational research workflows [17], [18].

Researchers have argued for over two decades that the complete software environment is a necessary condition for reproducing computational results [19]; however, the open sharing of such environments remains rare. Even when code is bug-free, compiles correctly, and is publicly shared, issues that are not resolved by code-sharing include (i) *code rot*, in which code becomes non-functional or its functionality changes as the underlying dependencies change over time (for example, an update to a data processing library which breaks backwards compatibility, or a modified implementation of an algorithm which changes experimental results), as well as (ii) *dependency hell*, in which configuring the software dependencies necessary to install or run code prevents successful execution [20]. This complex web of interdependencies is rarely described or documented in

published machine learning and computational science work [17], [9], [10].

**Methodological challenges** to reproducibility reflect the methods of the study, such as its procedure for model tuning or statistical evaluation. Existing work on reproducibility focuses on technical challenges, but methodological issues are at least as important. Methodological challenges include the use of biased model evaluation procedures [21], [22], the use of improperly-calibrated statistical tests for classifier comparison [23], or "large-scale hypothesis testing" where thousands of hypotheses or models are tested at once, despite the fact that most multiple testing corrections are not appropriate for such tasks [24]. A machine learning version of these errors is seen in massive unreported searches of the hyperparameter space, and in "random seed hacking" wherein the random number generator itself is systematically searched in order to make a target model's performance appear best or a baseline model worse [25]. Replication platforms can address methodological issues – working toward what [9] terms *inferential reproducibility*– by architecting platforms which support effective methodologies and adopt them by default, effectively nudging researchers to make sound choices.

**Data challenges** to reproducibility concern the availability of data itself. In some domains of big data research, making raw data available is more an issue of convention than a true barrier to reproducibility. However, educational data are governed by strict privacy regulations which protect the privacy of student records. Similar restrictions affect other big data domains, from the health sciences to computational nuclear physics [18]. As a result, researchers are often legally prohibited from making their data available. Efforts such as the Pittsburgh Science of Learning Center DataShop [26] and the HarvardX MOOC data sets [27] have attempted to address this problem in educational research by only releasing limited non-reidentiable data, but many analyses require the original, unprocessed data for a full replication. Restricted data sharing is one of the main factors (in our experience) hindering generalizability analysis in educational data mining: investigators are generally limited to one or two courses worth of data (e.g. the courses they instruct), and models are often overfit to these datasets.

## III. THE MOOC REPLICATION FRAMEWORK (MORF)

MORF's key design principles are: (i) ensure the complete end-to-end reproducibility of experiments, (ii) support high-quality replication studies and original research using multiple methods of analysis, (iii) provide access to a large and diverse dataset, (iv) leverage high-performance computing requiring minimal user input, and (v) ensure complete compliance with legal restrictions on data sharing. This section detail's MORF's architecture and how it encodes these principles. In Section IV, we discuss how this architecture

may apply to a much broader set of problems across many other domains.

### A. Platform Architecture

MORF consists of two main components: an open-source Python API for specifying the workflow of an experiment (the "MORF API"), and a Platform-as-a-Service (PaaS), which is a running instance of MORF's back-end infrastructure coupled with computational resources and a large MOOC dataset (the "MORF platform").

*1) MORF API and Controller Scripts:* The life-cycle of a complete end-to-end experiment, from raw data to results, is shown in Figure 1. First, a user creates and submits a configuration file to MORF, either using an HTTP request or using the `easy_submit()` MORF API function. This configuration file contains job metadata, including a pointer to an executable Docker image which encapsulates all code, software, and operating system dependencies for the users' experiment. The configuration file also points to a Python controller script that specifies the high-level experimental workflow, such as how model training and testing should occur and whether cross-validation or a holdout set should be used in a predictive modeling experiment. The use of controller scripts is a best practice for reproducible computational research [18], as it provides a single script to fully reproduce an experiment. An additional advantage is that MORF controller scripts are human-readable, providing a high-level overview of an experiment. An example of a controller script from the experiment in [28] is shown in Listing 1.

MORF uses the controller script to manage low-level data platform tasks, including (i) data wrangling (retrieving and archiving necessary data at each step of the experiment); (ii) Docker image setup and execution; and (iii) parallelization. Of particular note is (iii), as MORF is able to leverage parallelization without requiring any user input or code written for parallel execution. The controller script provides sufficient information about how MORF can execute parallelization, which can lead to speedups of 1-2 orders of magnitude when each of MORF's CPUs is occupied with a separate task (e.g. training models on each of the different MOOC courses available).

Listing 1: A MORF controller script

```
extract_session()
extract_holdout_session()
train_course(label_type = 'dropout')
test_course(label_type = 'dropout')
evaluate_course(label_type = 'dropout')
```

*2) Docker Containerization for Reproducibility:* Another key component of MORF's architecture is that it requires submissions as executable Docker containers in order to resolve many of the experimental challenges described
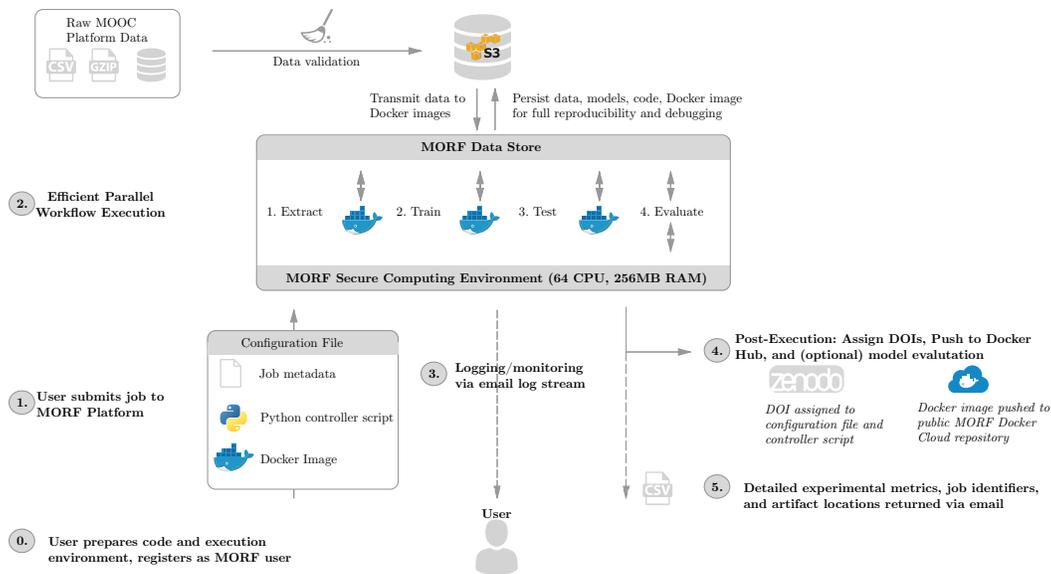
Figure 1: MORF platform architecture: this figure shows a high-level overview of a predictive modeling workflow in MORF. Solid lines indicate data transmission/archiving for reproducibility; dashed lines indicate messaging to MORF users.

in Section II-C. Docker containers can be thought of as lightweight virtual machines which fully encapsulate the code, software dependencies, and execution environment of an end-to-end experiment in a single file.

Docker containers were originally developed to resolve the technical issues described above in software development contexts [29], and are frequently used in both industrial software applications as well as computational and computer systems research [20], [30], [31]. Their use in data science applications is increasing, but the execution, publication, and sharing of pre-built Docker images as part of a research workflow is rare. While some existing platforms utilize containerization (e.g. Codalab[1]), this functionality is hidden from the user and containers are created "under the hood" by the platform. This limits users' ability to fully leverage containerization by building the complex, customized environments many machine learning experiments may require or to share these containers upon completion of an experiment. We are not aware of any data research platform which allow users to submit Docker images directly for execution.

A major advantage of Docker over simple code-sharing is that Docker containers fully reproduce the entire execution environment of the experiment, including code, software dependencies, and operating system libraries, exactly as this environment is configured at the time of an experiment. These containers are much more lightweight than a full virtual machine, but achieve the same level of reproducibility [29], [31].

When submitting a job for execution to the MORF platform, a user generates a Docker image containing the code, software, and operating system dependencies required to execute their experiment, and uploads the image to a public location (files located locally, HTTP, or in Amazon S3 are supported). The user provides the image's URL the configuration file submitted to MORF, and the image is fetched, checked, and executed according to the controller script. When an experiment completes error-free execution, MORF uploads the image to a public image repository on Docker Hub using a unique identifier. This makes implementations of every experiment on MORF immediately and publicly available for verification, extension, citation, or reuse.

As part of MORF, we are assembling an open-source library of ready-to-use Docker containers to replicate experiments conducted on MORF to serve as shared baseline implementations. These containers can be loaded with a single line of code, allowing the research community to replicate, fork, interrogate, modify, and extend the results of an experiment. For example, the experiment in [28] can be loaded by running `docker pull themorf/morf-public:fy2015-replication` in the terminal of any computer with Docker installed.

Building Docker containers requires only a single Dockerfile (akin to a makefile) which contains instructions for building the environment. This imposes minimal additional burden on researchers, but achieves a considerable increase in reproducibility over merely sharing code. Users can also generate Docker images by manually building a Docker environment (without a Dockerfile), or by using a tool such as Jupyter's `repo2docker`, which generates Docker images from public Github repositories[2].

---

[1]http://codalab.org/

[2]https://github.com/jupyter/repo2docker

MORF's combination of containerization and controller scripts allow the user to manage low-level experimental details (operating system and software dependencies, feature engineering methods, and statistical modeling) by constructing the Docker container which is submitted to MORF for execution, while MORF manages high-level implementation details (parallelization, data wrangling, caching of results) using the instructions provided in the controller script. Most importantly, the use of Docker containers ensures end-to-end reproducibility and enables sharing of the containerized experiment.

*3) Other Services and Platform Architecture:* While Docker containers allow jobs submitted to the MORF platform to use any combination of programming languages, software, or analytical tools that can be installed on a linux system, MORF itself is written in Python and the platform utilizes a variety of Amazon Web Services tools, including S3, Lambda, Simple Queueing Service (SQS), Simple Email Service (SES), and others. The platform is built on top of the Flask web framework. MORF utilizes Zenodo to assign unique DOIs to every experimental artifact (configuration files and controller scripts), and Docker Hub to create a public archive of every Docker image executed on the platform.

### B. Platform Use and Software API

MORF includes a simple, minimal software API which is used to write controller scripts that control job execution on the platform (see Listing 1 for an example). The Python API allows users to provide a simple execution "recipe" for the MORF platform to execute their experiment specifying the complete end-to-end pipeline from raw data to model evaluation: `extract` features from raw data; `train` and `test` machine learning models (predictive modeling experiments only), and `evaluate` the experimental results. For example, after extracting the desired features, a predictive modeling experiment could train individual models for every session of a course by using `train_session()` in their controller script; train one model per course using the data from all sessions by using `train_course()`; or train a single monolithic model using all data from every session of every course by using `train_all()`. An example of how MORF translates the various `train` functions into different experimental workflows, by mounting different data into the Docker environment, is shown in Figure 2. Note that the API functions used in the controller script are also used for MORF to manage many other low-level tasks, such as caching of intermediate results (e.g. trained models) and parallelization to ensure optimal performance. Similar functions exist for feature extraction, model testing, and evaluation, as well as for special cases (e.g. `forking` features extracted in a previous experiment for a new experiment).
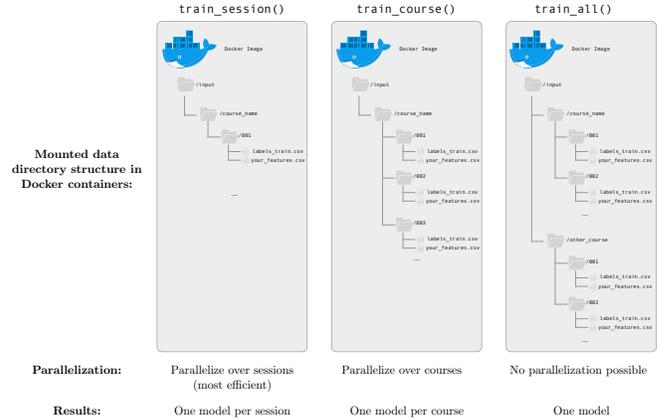
The MORF API is written in Python, is available on



Figure 2: Example of MOOC data mounted to Docker images in MORF based on the use of varying API functions in a predictive modeling experiment. This shows how users control workflows using the Python API, while having the flexibility to process and model raw data using any underlying software within their Docker images.

the Python Package Index (PyPi)[3], and is `pip`-installable. Development releases are available from Github.[4] It is important to note that while the MORF API is written in Python, this is simply the language that is required for controller scripts which MORF interprets to control job execution; users' code inside their Docker containers can be written using any language(s) or software that can be installed in the linux-based Docker container.

### C. Replication and Analysis Functionality

*1) Predictive Modeling:* A critical area of MOOC research to date has been the construction and analysis of predictive models of student success [8]. Such models have the potential to drive personalized learner supports or platform modalities [32], adaptive learning pathways [33], contribute to learning theory or support data understanding (such as about demographic differences in dropout or achievement) [4], or "early warning" systems designed to alert instructors of struggling students.

One primary thread of this research has been predicting whether a student is likely to dropout of, or fail to complete, a course. MORF currently supports two dropout prediction tasks: binary classification (dropout/no dropout) and regression (the week a student will drop out of a course, up to the final week). Predictive modeling experiments on MORF follow a standard end-to-end supervised learning workflow: feature extraction from raw data; model training; model testing; and model evaluation (whereby performance is analyzed or, optionally, evaluated using statistical tests). In order to perform a predictive modeling experiment using

---

[3]https://pypi.org/project/morf-api/
[4]https://github.com/educational-technology-collective/morf

MORF, users specify their Docker image should respond when executed in each of these "modes" with the requisite data mounted within the image. Predictive modeling experiments conducted on MORF to date include [28], [15].

*2) Production Rule Analysis:* A key design principle of MORF has been to support high-quality replication research, and the generation of other novel findings, even for researchers not interested in (or technically capable of) predictive modeling. As such, MORF supports a second mode of evaluation, known as *production rule* analysis. Production rules are if-then rules, coded in the form of "If a student who is <attribute> does <operator>, then <outcome>" [34]. MORF uses the JESS expert system language to code these production rules. In order to perform a production rule experiment using MORF, users provide a text file containing the production rule of interest. A set of attributes, operators, and outcomes are available in the MORF platform. Currently, the available attributes, operators, and outcomes are based on prior findings which we sought to replicate in the MORF platform, although future versions of the platform will allow users to specify their own definitions for each. Production rule experiments conducted on MORF to date include [34], [14].

### D. Platform Data

*1) Available MOOC Data:* To date, the privacy restrictions protecting MOOC learner data have hindered researcher access to large, diverse MOOC datasets. This has led to many studies evaluating only a small number of courses: a recent survey indicated that 70% of prior predictive modeling experiments in MOOCs used datasets consisting of five or fewer unique courses [8]. As discussed above, larger, more diverse and representative datasets are critical for a nascent field such as MOOC research, which is still developing consensus on many research questions.

The MORF platform seeks to resolve this issue, and makes a large, diverse dataset available for analysis of platform users. Currently, MORF contains the complete raw data exports from 209 sessions of 77 unique MOOCs offered on the Coursera platform, representing nearly 30 million unique interactions on courses offered across two of Coursera's four founding partner institutions (the University of Michigan and the University of Pennsylvania). Descriptive statistics for the data currently available in MORF are shown in Table I. The provision of this data for analysis to MORF users represents a considerable advancement of the field.

For all of the courses in MORF, the complete, raw exports consisting of all available course, user, and interaction data are available. In the case of the Coursera data currently made available by MORF, this consists of seven exports for every course [35]: multiple MySQL database dumps, CSV files, compressed clickstream server logs, and HTML files. This includes the complete clickstream server log for every user interaction with course content (example shown in Figure

Table I: Data available in the MORF Platform. This allows all experiments conducted on MORF to utilize massive, diverse datasets considerably larger than almost all other MOOC prediction research to date [8].

| Metric | Value |
|---|---|
| Total MOOC Sessions | 209 |
| Unique Courses | 77 |
| Active Students | 986,420 |
| Interactions | 29,416,369 |
| Discussion Forum Posts | 1,695,297 |
| Assignments | 231,066 |

3a), discussion forum posts and associated metadata (time, thread ID, user ID, upvotes/downvotes received, etc.; shown in Figure 3b), learner-provided responses to a course demographic survey, assignment submissions and grades with associated metadata, and course metadata (such as video titles, lengths and subtitle text; assignment information and due dates; etc.). Anonymizing such data, as many existing MOOC research solutions do, often obscures critical details, such as users' IP address (which can be used to estimate a users' location or socioeconomic status), the text of students' forum posts, and logfile entries which could be used to identify individual users. Using these rich datasets in raw and unredacted form presents an opportunity for researchers to investigate a highly diverse space of hypotheses, conduct rich feature engineering to build diverse predictive models, and replicate a wide variety of findings across many modes of inquiry.

*2) Execute-Against Access for Secure Big Data Research:* In order to provide direct and unredacted access to the rich, raw big data exports of MOOC platforms in a way that also accommodates the data regulations restricting the sharing of MOOC learner data (FERPA, IRB Common Rule, GDPR), we provide "execute-against" access to the platform data exports: users can run analyses against MORF's data, but cannot download or directly access the data. Instead, their experiments are executed against MORF's data within a sandboxed, networked-restricted computing environment. Then, a summary of the results are returned to the user. For predictive modeling jobs, these results are 8 measures of prediction performance (AUC, Cohen's $\kappa$, F1, etc.) on each course or session the job was executed on. For production rule analyses, statistical testing results are returned.

Most MOOCs are generated by a small number of platforms (e.g. Coursera, edX, FutureLearn), and all courses from a given platform use publicly-documented data schemas (e.g. [35]). Thus, users can develop experiments using their own data from a given platform – or even the public documentation – and then submit these experiments for MORF to execute against *any* other course from that platform. This enables MORF to provide an interface to its large data repository, without sharing the data itself, by utilizing the consistent and public schema of MOOC

{"key": "pageview","value": "https://www.coursera.org/lecture/39","username":
"                              ","timestamp": 1423101144446,"page_url": "https://
class.coursera.org/aidsfearandhope-001/lecture/39","client": "spark","session":
"8801600464-1422359379358","language": "en-US,en;q=0.5","from": "https://class.coursera.org/
aidsfearandhope-001/lecture/37","user_ip": "              ","user_agent": "Mozilla/5.0 (Windows NT 6.2;
WOW64; rv:35.0) Gecko/20100101 Firefox/35.0", "12": ["{\"height\":900,\"width\":1600}"], "13": [0], "30":
[1423101145091]}

{"key": "user.video.lecture.action","value": "{\"currentTime\":238.836,\"playbackRate\":1,\"paused\":true,
\"error\":null,\"networkState\":2,\"readyState\":3,\"lectureID\":11,\"eventTimestamp\":
1423100962836,\"initTimrestamp\":1423100112509,\"type\":\"pause\",\"prevTime\":238.836}","username":
"                              ","timestamp": 1423100963246,"page_url": "https://
class.coursera.org/aidsfearandhope-001/lecture/view?lecture_id=11,"client": "spark","session":
"8801600464-1422359379358","language": "en-US,en;q=0.5","from": "https://class.coursera.org/
aidsfearandhope-001/lecture/11","user_ip": "              ","user_agent": "Mozilla/5.0 (Windows NT 6.2;
WOW64; rv:35.0) Gecko/20100101 Firefox/35.0", "12": ["{\"height\":900,\"width\":1600}"], "13": [0], "30":
[1423100962836]}

(a)



(b)

Figure 3: Data types available for MOOC research in MORF. In existing solutions, these data sources are often highly redacted, anonymized, or simply unavailable. (a) Example clickstream log entries. (b) Example threaded forum posts. Personally-identifying information has been redacted.

datasets. These shared public data schemas also ensure that existing experiments in MORF can be replicated against new data (from the same MOOC platform) as it becomes available. Additionally, while only descriptive summary results are returned to users, MORF persists all extracted features, trained models, and predictions obtained in the course of an experiment for complete reproducibility, which also allows for trusted users to obtain access to intermediate experimental results.
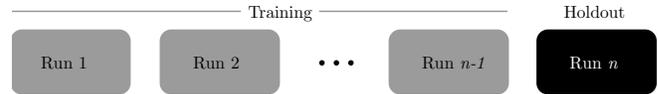


Figure 4: Default future-course prediction architecture used in MORF. This avoids potential inferential issues caused by using cross-validated performance to estimate model performance on future course sessions, by instead estimating that performance directly.

While the term "execute-against" access is novel, this framework of accepting analytical submissions which are executed against protected data has been used, for example, with healthcare data in the Critical Care Health Informatics Collaborative (CCHIC) [36] and with copyrighted music data in the Networked Environment for Music Analysis (NEMA) [37].

*E. MORF Resolves Existing Barriers to Replication in Educational Big Data Research*

MORF uses containerization to resolve many of the experimental challenges described in Section II-C. The Docker containers submitted to MORF fully encapsulate the code, software dependencies, and execution environment required of an experiment in a single file, ensuring end-to-end reproducibility and enabling sharing of the containerized experiment. These containers are automatically stored in a public repository and assigned unique, persistent identifiers. They can be loaded with a single line of code, allowing the research community to replicate, fork, interrogate, modify, and extend the results of any experiment in MORF.

In order to address concerns with methodological and inferential reproducibility, MORF provides sensible default procedures for many tasks, such as model evaluation, as well as simple statistical testing procedures. For example, MORF avoids the use of cross-validation for model evaluation: the prediction tasks to which most MOOC models aspire are prediction of *future* student performance (i.e., in an ongoing course where the true labels – such as whether a student will drop out – are unknown at the time of prediction). As such, using cross-validation within a MOOC session, when the outcome of interest is accuracy on a *future* MOOC session, provides an unrealistic and potentially misleading estimate of model performance. Prior work has demonstrated that cross-validation provides biased estimates of independent generalization performance [22], and in the MOOC domain, that cross-validation can produce biased estimates of classification performance on a future (unseen) course [28], [38]. Holding out a single session of every course requires a large data repository (multiple sessions of every MOOC). The rarity of such methods in prior work [39] is likely due to the scarcity of MOOC data, an issue MORF resolves. Adopting more effective model evaluation techniques by default requires no additional work for MORF users, and

<table>
<thead>
<tr><th></th><th>MORF</th><th>LearnSphere [26]*</th><th>CodaLab</th><th>CCHIC [36]</th></tr>
</thead>
<tbody>
<tr><td><strong>Data</strong></td><td></td><td></td><td></td><td></td></tr>
<tr><td>Raw (non-anonymized) data available</td><td>●</td><td>○</td><td>NA</td><td>○</td></tr>
<tr><td>Complete MOOC exports available</td><td>●</td><td>○</td><td>NA</td><td>NA</td></tr>
<tr><td>Total MOOC Sessions Available</td><td>209</td><td>151</td><td>NA</td><td>NA</td></tr>
<tr><td>Supports "execute-against" data</td><td>●</td><td></td><td></td><td>●</td></tr>
<tr><td><strong>Platform</strong></td><td></td><td></td><td></td><td></td></tr>
<tr><td>Supports custom code</td><td>●</td><td>●</td><td>●</td><td></td></tr>
<tr><td>Supports R, Python, Java, C++</td><td>●</td><td>●</td><td>●</td><td>○</td></tr>
<tr><td>Command-line submission/execution</td><td>●</td><td></td><td>●</td><td></td></tr>
<tr><td>Open-Source</td><td>●</td><td>●</td><td>●</td><td>●</td></tr>
<tr><td>GUI Available</td><td></td><td>●</td><td>●</td><td></td></tr>
<tr><td><strong>Reproducibility</strong></td><td></td><td></td><td></td><td></td></tr>
<tr><td>Uses containerization</td><td>●</td><td></td><td>●</td><td></td></tr>
<tr><td>Accepts submission as container</td><td>●</td><td></td><td></td><td></td></tr>
<tr><td>End-to-end execution in platform</td><td>●</td><td>○</td><td>●</td><td></td></tr>
<tr><td>Assigns DOI to each job</td><td>●</td><td></td><td></td><td></td></tr>
<tr><td>Analyses are forkable</td><td>●</td><td>○</td><td>●</td><td>○</td></tr>
<tr><td>Integrated w/public image repository</td><td>●</td><td></td><td></td><td></td></tr>
<tr><td><strong>Computation</strong></td><td></td><td></td><td></td><td></td></tr>
<tr><td>Provides computational resources</td><td>●</td><td>●</td><td>●</td><td></td></tr>
<tr><td>Native Parallelization Implemented</td><td>●</td><td></td><td></td><td></td></tr>
</tbody>
</table>

Table II: Comparison to existing solutions with similar goals, including those within the same domain (Learnsphere + Datastage), in a related domain (CCHIC), and domain-agnostic tools (Codalab). ●= Supported; ○= Partially Supported. *: statistics shown for use with DataStage MOOC data repository. †: users can provide their own additional worker nodes.

ensures that work produced on the MORF platform follows effective model evaluation procedures. MORF also provides support for statistical model evaluation to provide for more reliable inference than comparing machine learning experiments without statistical testing, as is commonly performed [39].

MORF achieves data reproducibility while also meeting data privacy restrictions by providing "execute-against" access to underlying data, described in Section III-D2.

## IV. DISCUSSION

### A. Implications for Big Data Research

In addition to jointly addressing several challenges to reproducible and replication research within the field of education, MORF's architecture, workflow, and initial research results have implications for the broader big data community. MORF addresses a set of problems faced by big data researchers across many domains. This includes experimental reproducibility as big data research in many fields uses increasingly complex computational models; methodological and inferential reproducibility as big data research enables problematic statistical practices such as massively multiple testing via testing thousands or millions of hypotheses in

a single experiment; and data reproducibility as available data become massively multimodal (many different formats), measure increasingly private or restricted aspects of users' behavior and identity, and cannot be easily anonymized.

MORF demonstrates that big data research in the face of such challenges is possible, and suggests a general blueprint for conducting a variety of research tasks while addressing these issues. MORF's API code is fully open-source, as is most of the code supporting the MORF platform (except where doing so presents a security risk). The MORF framework is domain-agnostic, and can support generic workflows for supervised learning and production rule analyses in any domain which works with complex, multiformat data which cannot be easily anonymized (e.g. sensitive medical data, copyrighted media, computational nuclear physics). As noted in Section III-D2 above, systems using execute-against access have been successfully deployed in the fields of healthcare and music information retrieval.

### B. Beyond Verification: The Benefits of Replication

Much prior work on reproducibility has focused on verifying that published results can be reproduced. However, end-to-end reproducible machine learning frameworks, such as MORF, provide benefits beyond mere verification, including:

**Gold standard benchmarking:** open replication platforms allow for the comparison of results which were previously not comparable, having been conducted on different data. The use of such benchmarking datasets has contributed to the rapid advance of fields such as computer vision (e.g. MNIST, IMAGENET), natural language processing (Penn Tree Bank, Brown corpus), and computational neuroscience (openFMRI). These datasets have been particularly impactful in fields where it is difficult or expensive to collect, label, or share data (as is the case with MOOC data, due to legal restrictions on sharing and access). These help to evaluate the state of the art by providing a common performance reference which is currently missing in many fields.

**Shared baseline implementations:** We noted above that variability in so-called "baseline" or reference implementations of prior work has contributed to concerns about reproducibility in the field [25]. By providing fully-executable versions of existing experiments, MORF ameliorates these issues, allowing for all future work to compare to the exact previous implementation of a baseline method.

**Forkability:** containerization produces a ready-made executable which fully encompasses the code and execution environment of an experiment. These can be readily shared and "forked" across the scientific community, much in the same way code is "forked" from a git repository. This allows researchers to build off of others' work by modifying part or all of an end-to-end pipeline (for example, by experimenting with different statistical algorithms but using the same feature set as a previous experiment) within the same software ecosystem.

**Generalizability analysis:** Each successive replication of an experiment provides information about its generalizability. Evaluating the generalizability of experiments has been a challenge in MOOC research to date, where studies conducted on single-course, restricted, and often homogenous datasets tend to dominate the literature. When containerized implementations are available, replicating these analyses on new data – even data which are not publicly available but share the schema of the original data – becomes as straightforward as running the containerized experiment against new data.

**Sensitivity Analysis:** This technique, used widely in Bayesian analysis, evaluates how changes to the underlying assumptions or hyperparameters affect experimental results. Such an evaluation can provide useful information about a model's robustness and potential to generalize to new data. Without being able to fully reproduce a model on the original data, sensitivity analyses of published results are not possible. In MORF, such analyses can be conducted by forking and modifying the containerized version of the original experiment, then re-executing it against the same data. This process also enables so-called *ablation analyses*, wherein individual components are removed from a model to observe their contribution to the results, as well as *slicing analyses*, where analysis of performance across different subgroups (e.g. demographic groups) is explored [40].

**Full Pipeline Evaluation:** Each stage of an end-to-end machine learning experiment (feature extraction, algorithm selection, model training, model evaluation) can be done in many different ways. Each stage also affects the others (for example, some algorithms might perform best with large feature spaces; others might perform poorly with many correlated features). However, current research usually evaluates only one or two components of this pipeline (e.g. training several algorithms and tuning their hyperparameters on a fixed feature set). Not only are the remaining stages often described in poor detail or not at all [9]; such work also leaves future researchers unable to evaluate the synergy between different aspects of the end-to-end pipeline in a published experiment (for example, exploring whether an algorithm's performance improves with a different feature set). MORF fully encapsulates this end-to-end pipeline for a given experiment and it makes it available for modification to any other researcher.

**Meta-Analysis:** While meta-analyses are common in fields with robust empirical research bases, such analyses have been less common in the field of big data, which has an emphasis on novelty. The open availability of executable machine learning experiments affords detailed meta-analyses by providing complete results of all modeling stages for meta-analysis. MORF has already been used for meta-analysis [14].

## V. Conclusion

Replication is important to all research, including big data research in education. We have constructed MORF, a system to facilitate such analysis. MORF has been used to conduct several replication studies [28], [15], [34], [14], each of which demonstrated that the initial findings did not entirely replicate. While we encourage interested institutions to partner with us to make their data available within MORF, users or institutions can also create their own platform instances by using the open-source code and API. We see immediate future opportunities to broaden our work by (a) reducing the need for researchers to understand environments such as Docker, by integrating front-end calls to this tool within development environments such as Jupyter; (b) providing light weight web front-ends for this system, allowing users to explore this large data through web visualizations; and (c) increasing the size and scope of the project by involving more institutions and a broader array of educational data (e.g. non-MOOC learning management system data).

## References

[1] M. Kloft, F. Stiehler, Z. Zheng, and N. Pinkwart, "Predicting MOOC dropout over weeks using machine learning methods," in *Proc. EMNLP 2014 Workshop on Analysis of Large Scale Social Interaction in MOOCs*.

[2] K. Veeramachaneni, U.-M. O'Reilly, and C. Taylor, "Towards feature engineering at scale for data from massive open online courses."

[3] D. Yang, T. Sinha, D. Adamson, and C. P. Rosé, "Turn on, tune in, drop out: Anticipating student dropouts in massive open online courses," in *Proc. 2013 NIPS Data-driven education workshop*, vol. 11, 2013, p. 14.

[4] R. F. Kizilcec and S. Halawa, "Attrition and achievement gaps in online learning," in *Proc. Learning@Scale*, 2015, pp. 57–66.

[5] C. G. Brinton and M. Chiang, "MOOC performance prediction via clickstream data and social learning networks," in *IEEE INFOCOM*, pp. 2299–2307.

[6] D. Yang, M. Wen, I. Howley, R. Kraut, and C. Rose, "Exploring the effect of confusion in discussion forums of massive open online courses," in *Proc. Learning@Scale*. ACM, pp. 121–130.

[7] M. L. Bote-Lorenzo and E. Gómez-Sánchez, "Predicting the decrease of engagement indicators in a MOOC," in *Proc. Seventh Intl. Learning Analytics & Knowledge Conference*. ACM, pp. 143–147.

[8] J. Gardner and C. Brooks, "Student success prediction in MOOCs," *UMUAI*, vol. 28, no. 2, pp. 127–203, 2018.

[9] O. E. Gundersen and S. Kjensmo, "State of the art: Reproducibility in artificial intelligence," in *AAAI-17 and the Twenty-Eighth Innovative Applications of AI Conf.*, 2017.

[10] B. K. Olorisade, P. Brereton, and P. Andras, "Reproducibility in machine Learning-Based studies: An example of text mining," in *Reproducibility in ML Workshop at the 34th Intl. Conf. on Machine Learning*, 2017.

[11] C. Collberg, T. Proebsting, G. Moraila, A. Shankaran, Z. Shi, and A. M. Warren, "Measuring reproducibility in computer systems research," Univ. Arizona Dept. of Comp. Sci., Tech. Rep., 2014.

[12] M. C. Makel and J. A. Plucker, "Facts are more important than novelty: Replication in the education sciences," *Educ. Res.*, vol. 43, no. 6, pp. 304–316, 2014.

[13] A. E. W. Johnson, T. J. Pollard, and R. G. Mark, "Reproducibility in critical care: a mortality prediction case study," in *Proc. 2nd Machine Learning for Healthcare Conference*, F. Doshi-Velez, J. Fackler, D. Kale, R. Ranganath, B. Wallace, and J. Wiens, Eds., vol. 68. PMLR, pp. 361–376.

[14] J. M. L. Andres, R. S. Baker, G. Siemens, D. Gašević, and S. Crossley, "Studying MOOC completion at scale using the MOOC replication framework," in *Proc. Intl. Conf. on Learning Analytics and Knowledge*, pp. 71–78.

[15] J. Gardner, C. Brooks, J. M. Andres, and R. Baker, "Replicating MOOC predictive models at scale," in *Proc. Fifth Annual Meeting of the ACM Conference on Learning@Scale*, 2018.

[16] V. Stodden and S. Miguez, "Best practices for computational science: Software infrastructure and environments for reproducible and extensible research," *Journal of Open Research Software*, vol. 2, no. 1, pp. 1–6, 2013.

[17] D. Donoho, "50 years of data science," in *Princeton NJ, Tukey Centennial Workshop*, 2015, pp. 1–41.

[18] J. Kitzes, D. Turek, and F. Deniz, *The Practice of Reproducible Research: Case Studies and Lessons from the Data-Intensive Sciences*. Univ of California Press.

[19] J. B. Buckheit and D. L. Donoho, "WaveLab and reproducible research," in *Wavelets and Statistics*, A. Antoniadis and G. Oppenheim, Eds. Springer New York, 1995, pp. 55–81.

[20] C. Boettiger, "An introduction to docker for reproducible research," *Oper. Syst. Rev.*, vol. 49, no. 1, pp. 71–79.

[21] G. C. Cawley and N. L. C. Talbot, "On over-fitting in model selection and subsequent selection bias in performance evaluation," *JMLR*, vol. 11, no. Jul, pp. 2079–2107, 2010.

[22] S. Varma and R. Simon, "Bias in error estimation when using cross-validation for model selection," *BMC Bioinformatics*, vol. 7, p. 91.

[23] T. G. Dietterich, "Approximate statistical tests for comparing supervised classification learning algorithms," *Neural Comput.*, vol. 10, no. 7, pp. 1895–1923.

[24] B. Efron and T. Hastie, *Computer Age Statistical Inference*. Cambridge University Press, 2016.

[25] P. Henderson, R. Islam, P. Bachman, J. Pineau, D. Precup, and D. Meger, "Deep reinforcement learning that matters," in *AAAI-18*, 2018, pp. 3207–3214.

[26] K. R. Koedinger, R. S. Baker, K. Cunningham, A. Skogsholm, B. Leber, and J. Stamper, "A data repository for the EDM community: The PSLC DataShop," *Handbook of educational data mining*, vol. 43, 2010.

[27] M. A. HarvardX, "HarvardX-MITx Person-Course academic year 2013 De-Identified dataset, version 2.0."

[28] J. Gardner, Y. Yang, R. Baker, and C. Brooks, "Enabling End-To-End machine learning replicability: A case study in educational data mining," in *Second Annual Workshop on Reproducibility in Machine Learning at the Thirty-fifth Intl. Conf. on Machine Learning*, 2018.

[29] D. Merkel, "Docker: Lightweight linux containers for consistent development and deployment," *Linux J.*, vol. 2014, no. 239.

[30] J. Cito, V. Ferme, and H. C. Gall, "Using docker containers to improve reproducibility in software and web engineering research," in *Web Eng.* Springer, Cham, pp. 609–612.

[31] D. M. Jacobsen and R. S. Canon, "Contain this, unleashing docker for hpc," *Proc. Cray User Group*, 2015.

[32] J. Whitehill, J. Williams, G. Lopez, C. Coleman, and J. Reich, "Beyond prediction: Toward automatic intervention to reduce MOOC student stopout," in *Proc. 8th Intl. Conf. on Educational Data Mining*, 2015, pp. 171–178.

[33] Z. A. Pardos, S. Tang, D. Davis, and C. V. Le, "Enabling Real-Time adaptivity in MOOCs with a personalized Next-Step recommendation framework," in *Proc. Learning@Scale*. ACM, pp. 23–32.

[34] J. M. L. Andres, R. S. Baker, G. Siemens, D. Gašević, and C. A. Spann, "Replicating 21 findings on student success in online learning," *Technology, Instruction, Cognition, and Learning*, pp. 313–333, 2016.

[35] Coursera, *Coursera Data Export Procedures*.

[36] S. Harris, S. Shi, D. Brealey, N. S. MacCallum, S. Denaxas, D. Perez-Suarez, A. Ercole, P. Watkinson, A. Jones, S. Ashworth, R. Beale, D. Young, S. Brett, and M. Singer, "Critical care health informatics collaborative (CCHIC): Data, tools and methods for reproducible research: A multi-centre UK intensive care database," *Int. J. Med. Inform.*, vol. 112, pp. 82–89.

[37] K. West, A. Kumar, A. Shirk, G. Zhu, J. S. Downie, A. Ehmann, and M. Bay, "The networked environment for music analysis (NEMA)," in *2010 6th IEEE World Congress on Services*, 2010, pp. 314–317.

[38] J. Whitehill, K. Mohan, D. Seaton, Y. Rosen, and D. Tingley, "Delving deeper into MOOC student dropout prediction."

[39] J. Gardner and C. Brooks, "Evaluating predictive models of student success: Closing the methodological gap," *The Journal of Learning Analytics*, vol. 5, no. 2, pp. 105–125, 2018.

[40] D. Sculley, J. Snoek, A. Wiltschko, and A. Rahimi, "Winner's curse? on pace, progress, and empirical rigor," in *ICLR 2018 Workshops*.