

MODELDB: Opportunities and Challenges in Managing Machine Learning Models

Manasi Vartak
MIT CSAIL

mvaratak@csail.mit.edu

Samuel Madden
MIT CSAIL

madden@csail.mit.edu

Abstract

*Machine learning applications have become ubiquitous in a variety of domains. Powering each of these ML applications are one or more machine learning models that are used to make key decisions or compute key quantities. The life-cycle of an ML model starts with data processing, going on to feature engineering, model experimentation, deployment, and maintenance. We call the process of tracking a model across all phases of its life-cycle as **model management**. In this paper, we discuss the need for model management systems, describe MODELDB, the first open-source model management system developed at MIT, and discuss the opportunities and challenges in managing models.*

1 Introduction

Machine learning has become ubiquitous in a variety of applications including voice assistants, self-driving cars, and recommendation systems. Each ML-based application employs one or more machine learning models that are used to make key decisions or compute key quantities such as recognizing spoken words, detecting a pedestrian on the road, and identifying the best products to recommend. Models are to ML-based applications what databases are to stateful web-applications; they are key for the correct functioning of these applications. Consequently, just like we have database management systems (or DBMSs) to manage state in applications, we find the need for centralized systems that manage models, i.e., *model management systems*.

To understand the requirements of a model management system, we begin with a brief overview of the life-cycle of a machine learning model. The life-cycle of an ML model can be considered to consist of five phases, namely: (1) *Data Preparation* - Obtaining the training and test data to develop a model; (2) *Feature Engineering* - Identifying or creating the appropriate descriptors from the input data (i.e., features) to be used by the model; (3) *Experimentation* - Experimenting with different models on the training and test data and choosing the best; (4) *Deployment* - Deploying the chosen model in a live system; (5) *Maintenance* - Monitoring the live model performance and updating the model as needed. At every phase of the model life-cycle and during each iteration, certain pieces of metadata are essential for ensuring the key requirements at that phase. For example, for Phases 1-2, the key requirement for a model in that phase is reproducibility, i.e., enough data about the data preparation and feature engineering must be tracked so that the said features can be re-created exactly. Similarly, for Phase 3, the key requirement, along with reproducibility, is experiment tracking so that a data scientists can choose the best model. For Phases 4-5, the key requirements deal with making model-related data widely accessible and ensuring that key metrics for the models are collected regularly.

Copyright 0000 IEEE. Personal use of this material is permitted. However, permission to reprint/republish this material for advertising or promotional purposes or for creating new collective works for resale or redistribution to servers or lists, or to reuse any copyrighted component of this work in other works must be obtained from the IEEE.

Bulletin of the IEEE Computer Society Technical Committee on Data Engineering

We define a model management system as one that follows a model (whether deployed or not) throughout these five phases of its life-cycle and captures relevant metadata at each step. For instance, since the key requirement during the Experimentation phase is to enable the data scientists to choose the best model, metadata captured in this phases includes items such as performance metrics obtained by the model, hyperparameter values used during training, etc. In contrast, for a model in the Deployment phase, metadata might include the version of the model, where it is deployed, and how to query it. Although one can imagine a model management system that stores models and supports prediction serving operations, the heterogeneity in models, as well as their hardware and software requirements make such a solution sub-optimal. Therefore, we take the view that model management systems are best suited to store key metadata about models throughout their life-cycle. Thus, **we define a model management system as a system that tracks metadata about models through the five phases of their life-cycle.**

Given the rapid proliferation of machine learning applications, systems have been proposed in academia as well as industry to address different aspects of the model management problem. In this paper, we focus on MODELDB, the first open-source system we developed at MIT for model management. Other academic systems that seek to address similar problems include the ModelHub system [17] (to explore deep learning architectures and efficiently store network weights), ProvDB [18] (to manage metadata collected via collaborative data science), Ground [13] (to provide a common framework for tracking data origin and use via generic abstractions), and the work on model selection management systems by Kumar et. al. [14].

One of the earliest commercial model management system was the SAS Model Manager which tracked models built and deployed within the SAS platform [23]. Today, most proprietary ML platforms such as the Michelangelo platform at Uber [8] and FBLearner at Facebook [7] include a model management or model repository component. Similarly, vast majority of data science teams end up re-building a model management system to suit their needs. In [24], Sculley et. al. elegantly present the challenges with building and productionizing at Google and highlight the need to manage “pipeline jungles” and model configurations.

Model management systems are also closely related to workflow and experiment management systems such as Kepler [15], Taverna workbench [28], Galaxy [26], VisTrails [1, 3, 2] as well as recent workflow engines tailored for data processing such as Apache Airflow [10] and Luigi [25]. While workflow systems can address some of the model management needs during the first three phases of the model life-cycle, these systems require extensions to support the deployment and maintenance phases of the ML model life-cycle.

The rest of the paper is organized as follows. In Section 2, we describe the motivation behind model management systems; in Section 3, we describe the challenges faced in building model management systems; in Section 4, we describe the MODELDB system developed at MIT. We conclude the paper in Section 5 with a discussion of how we see the need for model management evolving in the future.

2 Need for Model Management

To understand the need for model management, we studied modeling workflows across many companies and research labs ranging in size and complexity of their machine learning applications. Across the different phases of the model life-cycle, the need for model management was apparent in three key areas: managing modeling experiments, enabling reproducibility and collaboration, and the need for governance.

The first and primary area where the need for model management is evident is during the Experimentation phase of the model life-cycle. The empirical nature of machine learning model development means that data scientists and ML developers experiment with hundreds of models before identifying one that meets some acceptance criteria. This is particularly true when the data scientists performs hyperparameter optimization on a model. Data about previously built ML workflows is necessary to inform the next set of ML workflows and experiments, and also simply to identify the best model produced so far. Consequently, tracking these experiments becomes of paramount importance. The absence of an experiment tracking system leads to models and experi-

ments being *lost* and valuable time and resources being spent in reproducing the models. For example, one ML developer at a large tech company related how she had spent over a week just re-running a modeling experiment another employee had previously conducted since the experimental setup and results were not recorded.

A second related need for model management becomes evident when a previously built model needs to be reproduced from scratch. For example, perhaps a new version of a model produces erroneous results and therefore the model must be reverted to an older version. If the previous model object has been removed, re-creating it requires accurate information about what data was used, how it was processed, libraries and versions used, and details of how the model was trained (including random seeds). A similar need becomes evident when there is a discrepancy between an offline and live model or when a model has to be updated with new data. Similarly, when data science teams want to collaborate on a particular model or build on top of each other's work (e.g., use the results of one model as input to another), the lack of a centralized repository of models hampers sharing of information about what models exist and how to integrate a given models in a product or business process.

Another requirement for model management arises from the need to provide governance around what models are being used to make automated decisions. For non-regulated industries, this governance may only be required to provide the business insight into what models are being used in a product or business. However, for many regulated industries, particularly in banking and healthcare, government regulations require that any models used to make automated decisions be documented and available for audits. Moreover, government legislation now being implemented (e.g., the GDPR regulations in the European Union [20]) require that companies be able to explain any decisions made without human intervention, making it essential to record all models used across a business and provide complete provenance information about how the model was generated and used.

3 Challenges in Model Management

As defined above, model management covers the entire life-cycle of models starting with data-processing, feature engineering, model experimentation, deployment, and maintenance. Across all the above stages, *diversity* and *heterogeneity* in ML environments and frameworks are the key challenges in consistently capturing metadata across the modeling lifecycle.

For example, during data processing (or feature engineering), we want to track the transformations applied to data so that they can be accurately reproduced later. Since feature processing may be done in very different languages and compute environments (e.g., Spark, Teradata, HBase), ensuring that data processing can be accurately recorded across these environments is challenging. Moreover, even within a single environment, ensuring high coverage for all the ways in which a data transformation may be applied (e.g., applying a one-hot-encoding operation in pandas) is challenging unless this functionality has been built into the system from the ground-up.

Similarly, during the experimentation phase, there is a large diversity in machine learning environments and libraries in different languages (e.g., scikit-learn, Tensorflow, PyTorch in Python; R; H2O framework in Java). Each library has a unique way of defining models (e.g., graphs in Tensorflow vs. plain objects in scikit-learn) and their associated attributes. Consequently, one representation cannot capture models built across all frameworks. The diversity in machine learning frameworks also translates to diversity in deployment methods for each model. For example, while TF-serving is a popular means to serve Tensorflow models, a Flask-based deployment method is most popular for scikit-learn models. Capturing the relevant deployment properties for the above methods therefore also becomes a challenge.

Once a model is deployed, while some properties are common to all models (e.g., latency of predictions, number of requests), many properties are application dependent (e.g., application-specific metrics to be monitored) and might require input from disparate systems (e.g., prediction storage systems). As a result, one generic solution is unlikely to meet the requirements for all applications and will instead requires extensible code.

In addition to the challenges introduced due to the diversity of ML environments, two overarching require-

```

# version 61
#   Drop fireplacecnt and fireplaceflag, following Jayaraman:
#   https://www.kaggle.com/valadi/xgb-w-o-outliers-lgb-with-outliers-combo-tune5

# version 60
#   Try BASELINE_PRED=0.0115, since that's the actual baseline from
#   https://www.kaggle.com/aharless/oleg-s-original-better-baseline

# version 59
#   Looks like 0.0056 is the optimum BASELINE_WEIGHT

# versions 57, 58
#   Playing with BASELINE_WEIGHT parameter:
#   3 values will determine quadratic approximation of optimum

...

# version 49
#   My latest quadratic approximation is concave, so I'm just taking
#   a shot in the dark with lgb_weight=.3

# version 45
#   Increase lgb_weight to 0.25 based on new quadratic approximation.
#   Based on scores for versions 41, 43, and 44, the optimum is 0.261
#   if I've done the calculations right.
#   I'm being conservative and only going 2/3 of the way there.
#   (FWIW my best guess is that even this will get a worse score,
#   but you gotta pay some attention to the math.)

# version 44
#   Increase lgb_weight to 0.23, per Nikunj's suggestion, even though
#   my quadratic approximation said I was already at the optimum

```

Figure 1: Model Versioning comments by Kaggle competitor

ments emerge from requirements for model management: first, data scientists want a model management solution that minimizes the developer intervention and requires minimal changes to the current developer workflow. For example, many data scientists are (understandably) unwilling to choose a new ML environment or significantly change their modeling workflow to account for model management. Second, data scientists want a *vendor-neutral* model management system so that users are not tied into one particular provider or framework and can instead focus on their modeling task.

4 MODELDB

Building an ML model for real-world applications is an iterative process. Data scientists and ML developers experiment with tens to hundreds of models before identifying one that meets some acceptance criteria on model performance. For example, top competitors in the Kaggle competition to predict Zillow Home prices [5] made more than 250 submissions (and therefore built at least as many models), while those in the Toxic Comment classification competition [4] made over 400 submissions. As an example of actual experimentation performed during model building, Listing 1 reproduces code comments by an expert Kaggle competitor (ranked 500) written in order to track models built for the Zillow Price Prediction Challenge.

As we can see from this listing, a data scientist typically tests a large number of model versions before identifying the best one. Moreover, although data scientists (and data science teams) build many tens to hundreds of models when developing an ML application, they currently have no way to keep track of all the models they have built. Consequently, insights are lost, models cannot be reproduced, collaboration becomes challenging, and model governance becomes challenging.

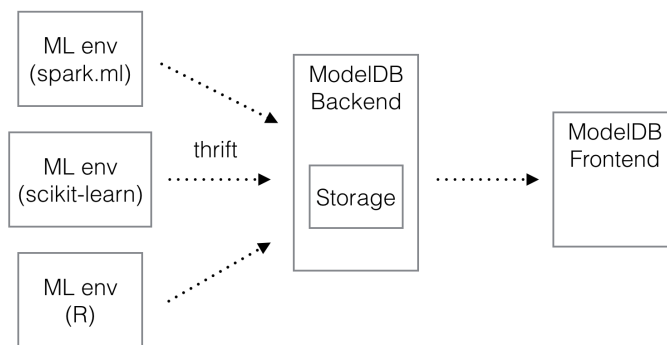


Figure 2: ModelDB Architecture

To address these challenges, we developed a system at MIT called MODELDB [27]. MODELDB is the first open-source machine learning model management system and currently focuses on tracking models during the experimentation phase. MODELDB automatically tracks models as they are built, records provenance information for each step in the pipeline used to generate the model, stores this data in a standard format, makes it available for querying via an API and a visual interface. In the following sections, we describe the architecture of MODELDB, how data may be logged into the system and different ways of querying the data. We finish with a brief discussion of current work on extending MODELDB.

4.1 MODELDB Architecture

Figure 2 shows the high-level architecture of our system. MODELDB consists of three key components: native client libraries for different machine learning environments, a backend that stores model data, and a web-based visualization interface. Client libraries are responsible for automatically extracting models and pipelines from code and passing them to the MODELDB backend. MODELDB client libraries are currently available for scikit-learn, spark.ml, and along with a *Light* Python API that can be used in any Python-based machine learning environment. This means that ML developers can continue to build models and perform experimentation these environments while the native libraries passively capture workflow information. The MODELDB backend exposes a thrift¹ interface to allow clients in different languages to communicate with the MODELDB backend. The backend can use a variety of storage systems to store the model metadata. The third component of MODELDB, the visual interface, provides an easy-to-navigate layer on top of the backend storage system that permits visual exploration and analyses of model data.

4.2 Client Libraries

Many existing workflow management programs (e.g. VisTrails [3]) require that the user create a workflow in advance, usually by means of a GUI. However, the data scientists we interviewed overwhelmingly concurred that GUIs restricted their flexibility in defining pipelines and that it was difficult to specify pipelines beforehand because they changed constantly. Moreover, we found that data scientists were unwilling to change their preferred ML environment for a workflow management system. Therefore, our primary design constraint while creating the MODELDB client libraries was to minimize any changes the data scientist would need to make both to code and the existing modeling process. To meet this constraint, we chose to make model logging accessible directly through code (as opposed to a GUI) and to build native logging libraries for different ML environment.

¹<https://thrift.apache.org/>

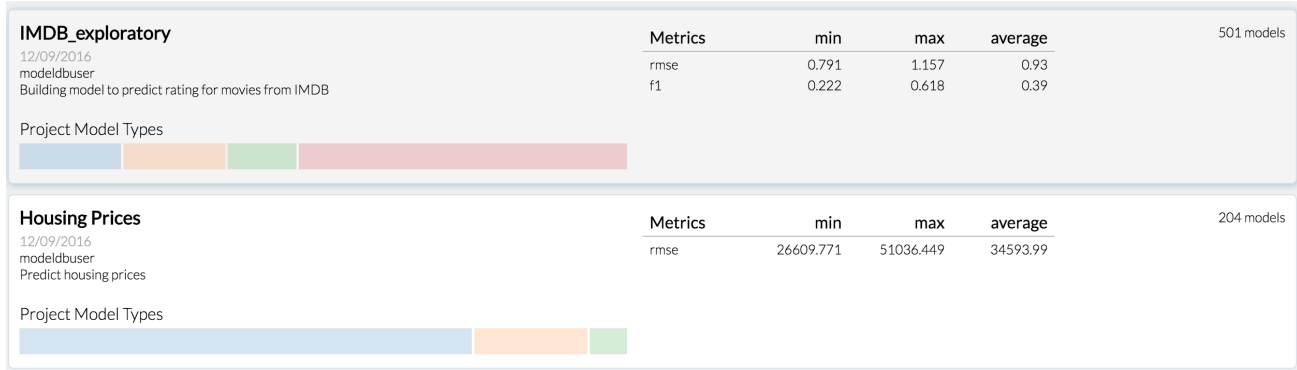


Figure 3: Projects Summary View

The spark.ml and scikit-learn libraries are architected such that data scientists can use the environments for analysis exactly as they normally would and the library transparently and automatically logs the relevant data to the backend.

4.3 MODELDB Frontend

MODELDB captures a large amount of metadata about models. To support easy access to this data, MODELDB provides a visual interface. We provide three key views for exploring the data stored in MODELDB. A user starts with the projects summary page (Fig. 3) that provides a high level overview of all the projects in the system. The user can then click on a particular project to see the workflows for that project. We present models via two key views; the first view presents two visualizations, a timeline visualization (Fig. 4) showing the evolution of the model metrics over time as well as a custom visualization interface for model meta-analyses. The second view is a tabular view of all the models in a particular project (Fig. 5) along with interactions such as filtering, sorting, grouping, and search. From any of the above interfaces, the ML developer can drill-down into a single model to visualize the model pipeline that was automatically inferred by the client library.

4.4 MODELDB Adoption and Future Work

We officially released MODELDB as an open-source model management system in Feb. 2017 and since then there has been a large amount of interest and adoption of MODELDB. Specifically, over the last year, our GitHub repository [11] has garnered > 500 stars, has been cloned over a thousand times, and has been forked >100 times. MODELDB has been tested at multiple small and large companies, and is deployed in various settings. MODELDB has also served as inspiration for other model management systems such as [19].

Based on the feedback from current MODELDB users and broader data science needs, we are currently expanding MODELDB to capture metadata not only from the Experimentation phase, but also the Deployment and Monitoring phases.

5 Evolution of Model Management Systems

Current model management systems have focused on keeping track of modeling experiments and, to some extent, keeping track of deployed model versions. However, as machine learning models proliferate into every key business process and product, we posit that the task of managing the life-cycle of models will become as key as managing the life-cycle of code. Just as version control systems such as SVN and Git made source code development and collaboration more robust, we envision that model management systems will serve a similar



Figure 4: Timeline Visualization

IDs	DataFrame	Specifications	Metrics	Misc.
Model ID: 22 Experiment Run ID: 1 Experiment ID: 1	DataFrame ID: 30	Type: LinearRegression Hyperparameters	rmse: 0.881	Notes: test annotation [+] Model Filepath: 2016-12-09 ... Timestamp: 2016-12-09 18:5... See Mo
Model ID: 29 Experiment Run ID: 1 Experiment ID: 1	DataFrame ID: 35	Type: LinearRegression Hyperparameters	rmse: 0.849	Notes: this model is funky [+] Model Filepath: 2016-12-09 ... Timestamp: 2016-12-09 18:5... See Mo
Model ID: 30 Experiment Run ID: 1 Experiment ID: 1	DataFrame ID: 36	Type: LinearRegression Hyperparameters	rmse: 0.873	Notes: feature1, feature2, fe... [+] Model Filepath: 2016-12-09 ... Timestamp: 2016-12-09 18:5... See Mo
Model ID: 31 Experiment Run ID: 1 Experiment ID: 1	DataFrame ID: 37	Type: LinearRegression Hyperparameters	rmse: 0.942	Notes: [+] Model Filepath: 2016-12-09 ... Timestamp: 2016-12-09 18:5... See Mo

Figure 5: Models View

purpose in the future. We imagine model management systems to become the system of record for all models and model-related information. Consequently, we expect model management systems to evolve in the following directions.

Model Data Pipelines. For many machine learning applications, the ultimate performance of the machine learning model depends on the features or data attributes used by the model. As a result, when a model is to be reproduced, it requires accurate records of the data and transformations used to produce features. We expect model management systems to evolve to accurately record data versions (train and test) as well data transformations that are used to generate features. This metadata may come from a separate data processing or workflow system as opposed to being generated by the model management system, however, the model management system would track all the metadata required to create the model end-to-end.

Model Interoperability. One of the key challenges in model management described before is the diversity in ML models and frameworks. While some frameworks support rapid development, others might be more suitable for deployment in production settings. Given the rapid proliferation of ML frameworks, we expect model management systems to support, if not provide, a level of interoperability between different ML environments and frameworks (e.g., PMML [12] and ONNX [21] provide a good start). This will enable data scientists to pick and choose the best framework for each phase of the model life-cycle.

Model Testing. As more decisions and business logic get delegated to machine learning models, the importance of testing models (similar to code testing) will increase and will become a key part of managing the model life-cycle. For instance, defining unit tests for models and their input data will become commonplace. Similarly, we expect integration tests with model outputs to become more prevalent as model predictions get used as input for other models. Finally, we note that as adversarial attacks on models increase, testing of models and edge cases will become key, requiring the development of new techniques to prevent adversarial attacks (e.g., [?, 9]).

Model Monitoring. While model testing takes place before a model is deployed, model monitoring takes place once the model is deployed in a live system. Model monitoring today is largely limited to monitoring system-level metrics such as the number of prediction requests, latency, and compute usage. However, we are already seeing the need for data-level monitoring of models (e.g., as noted in [24]) to ensure that the offline and live data fed to a model is similar. We expect model management systems to encompass model monitoring modules to ensure continued model health, triggering alerts and actions as appropriate.

Model Interpretability and Fairness. As models are used for automated decision making in regulated industries and increasingly used by non-technical users, explaining the results of models will become a key aspect of the management of deployed models (as evidenced by the rich research on interpretability [16, 6, 22]). We view a model management system as the system of record for all models and, therefore, the logical gateway for model interpretability and understanding. In the future, we therefore expect model management systems to expose interpretability functionality for every recorded model.

6 Conclusion

In this paper, we discussed the need for model management systems that track model metadata throughout the model life-cycle. We described MODELDB, an open-source model management system developed at MIT that focuses on the Experimentation phase of the model life-cycle and is being extended to Deployment and Maintenance. Finally, we described the challenges and opportunities in model management in the future.

References

- [1] Louis Bavoil, Steven P Callahan, Patricia J Crossno, Juliana Freire, Carlos E Scheidegger, Cláudio T Silva, and Huy T Vo. Vistrails: Enabling interactive multiple-view visualizations. In *Visualization, 2005. VIS 05. IEEE*, pages 135–142. IEEE, 2005.

- [2] Steven P Callahan, Juliana Freire, Emanuele Santos, Carlos E Scheidegger, Claudio T Silva, and Huy T Vo. Managing the evolution of dataflows with vistrails. In *Data Engineering Workshops, 2006. Proceedings. 22nd International Conference on*, pages 71–71. IEEE, 2006.
- [3] Steven P Callahan, Juliana Freire, Emanuele Santos, Carlos E Scheidegger, Cláudio T Silva, and Huy T Vo. Vistrails: visualization meets data management. In *Proceedings of the 2006 ACM SIGMOD international conference on Management of data*, pages 745–747. ACM, 2006.
- [4] Kaggle competition. Toxic comment classification challenge. <https://www.kaggle.com/c/jigsaw-toxic-comment-classification-challenge>.
- [5] Kaggle Competition. Zillow prize: Zillows home value prediction (zestimate). <https://www.kaggle.com/c/zillow-prize-1>.
- [6] Been Doshi-Velez, Finale; Kim. Towards a rigorous science of interpretable machine learning. In *eprint arXiv:1702.08608*, 2017.
- [7] Facebook Engineering. Introducing fblearner flow: Facebook’s ai backbone. <https://code.facebook.com/posts/1072626246134461/introducing-fblearner-flow-facebook-s-ai-backbone/>.
- [8] Uber Engineering. Meet michelangelo: Uber’s machine learning platform. <https://eng.uber.com/michelangelo/>.
- [9] Reuben Feinman, Ryan R Curtin, Saurabh Shintre, and Andrew B Gardner. Detecting adversarial samples from artifacts. *arXiv preprint*, 2017.
- [10] Apache Software Foundation. Airflow.
- [11] MIT DB Group. Modeldb. <https://github.com/mitdbg/modeldb>, 2017.
- [12] Alex Guazzelli, Wen-Ching Lin, and Tridivesh Jena. *PMML in Action: Unleashing the Power of Open Standards for Data Mining and Predictive Analytics*. CreateSpace, Paramount, CA, 2010.
- [13] Joseph M Hellerstein, Vikram Sreekanti, Joseph E Gonzalez, James Dalton, Akon Dey, Sreyashi Nag, Krishna Ramachandran, Sudhanshu Arora, Arka Bhattacharyya, Shirshanka Das, et al. Ground: A data context service.
- [14] Arun Kumar, Robert McCann, Jeffrey Naughton, and Jignesh M. Patel. Model selection management systems: The next frontier of advanced analytics. *SIGMOD Rec.*, 44(4):17–22, May 2016.
- [15] Bertram Ludäscher, Ilkay Altintas, Chad Berkley, Dan Higgins, Efrat Jaeger, Matthew Jones, Edward A Lee, Jing Tao, and Yang Zhao. Scientific workflow management and the kepler system. *Concurrency and Computation: Practice and Experience*, 18(10):1039–1065, 2006.
- [16] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In *Advances in Neural Information Processing Systems 30*, pages 4768–4777. Curran Associates, Inc., 2017.
- [17] H. Miao, A. Li, L. S. Davis, and A. Deshpande. Modelhub: Deep learning lifecycle management. In *2017 IEEE 33rd International Conference on Data Engineering (ICDE)*, pages 1393–1394, April 2017.
- [18] Hui Miao, Amit Chavan, and Amol Deshpande. ProvdB: Lifecycle management of collaborative analysis workflows. 2017.
- [19] MLFlow. Mlflow. <https://github.com/mlflow/mlflow>.
- [20] Official Journal of the European Union. General data protection regulation. <https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=celex%3A32016R0679>.
- [21] ONNX. Onnx: Open neural network exchange. <https://github.com/onnx/onnx>.
- [22] Maithra Raghu, Justin Gilmer, Jason Yosinski, and Jascha Sohl-Dickstein. Svcca: Singular vector canonical correlation analysis for deep learning dynamics and interpretability. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 6078–6087. Curran Associates, Inc., 2017.
- [23] SAS. Sas model manager. https://www.sas.com/en_us/software/model-manager.html.

- [24] D Sculley, Todd Phillips, Dietmar Ebner, Vinay Chaudhary, and Michael Young. Machine learning: The high-interest credit card of technical debt.
- [25] Spotify. Luigi. <https://github.com/spotify/luigi>.
- [26] The Galaxy Team and Community. The galaxy project. <https://galaxyproject.org/>.
- [27] Manasi Vartak, Harihar Subramanyam, Wei-En Lee, Srinidhi Viswanathan, Saadiyah Husnoo, Samuel Madden, and Matei Zaharia. Modeldb: A system for machine learning model management. In *Proceedings of the Workshop on Human-In-the-Loop Data Analytics*, HILDA '16, pages 14:1–14:3, New York, NY, USA, 2016. ACM.
- [28] Katherine Wolstencroft, Robert Haines, Donal Fellows, Alan Williams, David Withers, Stuart Owen, Stian Soiland-Reyes, Ian Dunlop, Aleksandra Nenadic, Paul Fisher, et al. The taverna workflow suite: designing and executing workflows of web services on the desktop, web or in the cloud. *Nucleic acids research*, 41(W1):W557–W561, 2013.