

# GNG: Data-Centric AI Engineering

Arumoy Shome

[2022-08-30 Tue]

This document presents my research portfolio for the GNG meeting. The document presents an introduction to the research field of Software Engineering for Machine Learning (SE4ML), the intended research field for this thesis. An overview of the state-of-the-art is presented next followed by the research topic of choice. The document concludes with contributions that I have made so far both scientific & supplemental. I also list ideas that were unsuccessful along with ideas we wish to pursue in the future.

## 1 Background & Motivation

Machine Learning and database systems are well established research fields. Since ML models & datasets go hand-in-hand, their intersection has been well studied. However, with the growing adoption of ML beyond an academic setting, a need for engineering best practices is felt. Although decades of research exists for traditional rule-based software systems, the introduction of Software Engineering for ML systems has been a fairly recent event in the research community.

Figure 1 presents an overview of the Machine Learning Lifecycle. Data analysis is a critical and dominant stage of the machine learning lifecycle. Once the data is collected, most of the work goes into studying and wrangling the data to make it fit for training. A highly experimental phase follows where a model is selected and tuned for optimal performance. The final model is then productionised and monitored constantly to detect data drifts and drop in performance [9, 5, 2, 10].

The *data science* stage is a typical ML workflow which tends to dominate blogposts shared on the internet and also what we commonly see in an academic setting. However as seen in 2, the ML workflow tends to be a very small part of a much larger system.

# ML Lifecycle

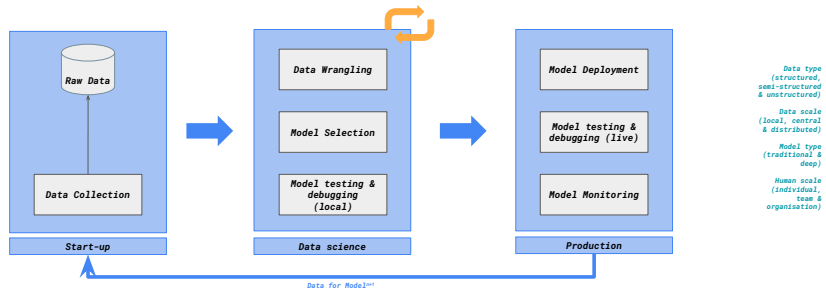
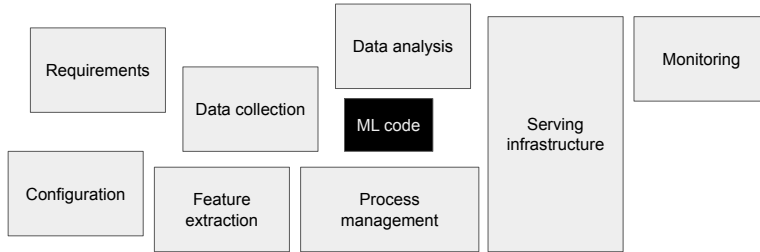


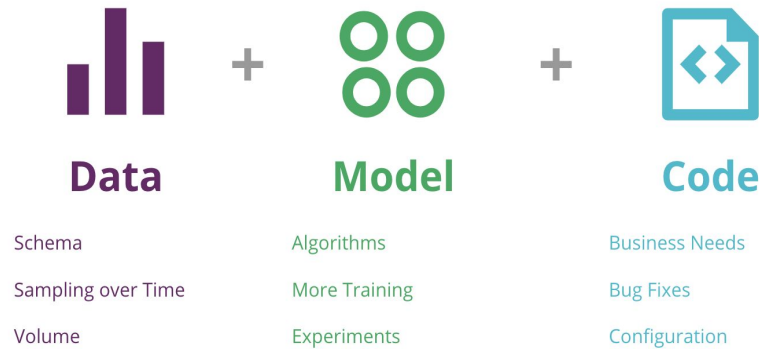
Figure 1: ML Lifecycle



D. Sculley et al., Hidden technical debt in machine learning systems

3

Figure 2: Holistic View of ML Systems



D. Sato et al., [Continuous Delivery for Machine Learning](#)

4

Figure 3: ML Pillars of Change

When compared to traditional software, the feedback loop of a machine learning system is longer. While traditional software primarily experiences change in *code*, a machine learning system matures through changes in *data*, *model* & *code* [9]. Given the highly tangled nature of machine learning systems, a change in any of the stages of the lifecycle triggers a ripple effect throughout the entire pipeline [10]. Testing such changes also becomes challenging since all three components need to be tested. Besides the traditional test suites, a full training-testing cycle is required which incurs time, resource and financial costs. The surrounding infrastructure of a machine learning pipeline becomes increasingly complex as we move towards a productionised model. Thus catching potential problems in the early, upstream phase of data analysis becomes extremely valuable as fixes are faster, easier and cheaper to implement.

AI has had a significant impact on the technology sector due to the presence of large quantities of unbiased data [7]. But AI’s true potential lies in its application in critical sectors such as healthcare, wildlife preservation, autonomous driving, and criminal justice system [3]. Such high-risk domains almost never have an existing dataset and require practitioners to collect data. Once the data is collected, it is often small and highly biased. While AI research is primarily dominated by model advancements, this new breed of *high-stakes AI* supports the need for a more data-centric approach to AI [6, 8, 11].

ML tends to be very data-centric in nature and majority of the work

involves working with the data. When we put a software engineering lens on ML problems, **data becomes equivalent to code**. While software engineers have several tools & techniques to aid them in their day-to-day lives, the same cannot be said for ML practitioners.

## 2 Scientific Contributions

This section presents the scientific contributions that we have made so far. I also include the ideas which we are currently working on, ideas that we wish to pursue in the future & ideas which were unsuccessful.

### 2.1 DONE shome2021data

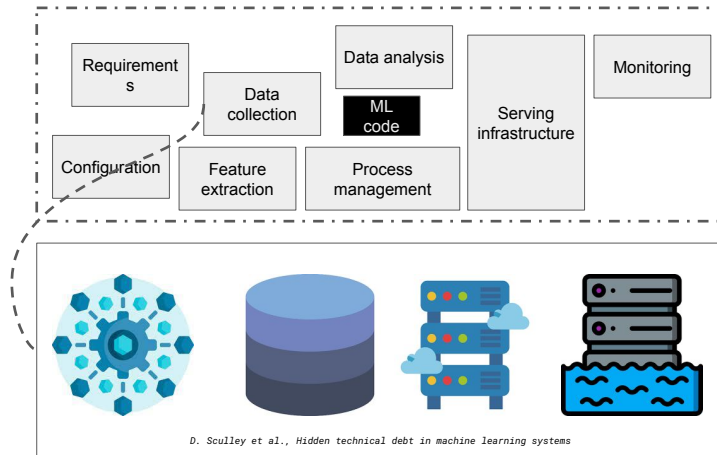
Since the study of software systems with machine learning components is a fairly young discipline, resources are lacking to aid practitioners in their day-to-day activities. The highly data-driven nature of machine learning makes data equivalent to code in traditional software. The notion of code smells is critical in software engineering to identify early indications of potential bugs, sources of technical debt and weak design choices. Code smells have existed for over 30 years. A large body of scientific work has catalogued the different smells, the context in which they occur and their potential side-effects. To the best of our knowledge, such a catalogue however did not exist for data science.

This study identified the recurrent data quality issues in public datasets. Analogous to code smells, we introduced a novel catalogue of data smells that can be used to indicate early signs of problems or technical debt in machine learning systems. To understand the prevalence of data quality issues in datasets, we analysed 25 public datasets and identified 14 data smells.

The paper was accepted at The 1<sup>st</sup> International Conference on AI Engineering '22, which was co-located with ICSE '22.

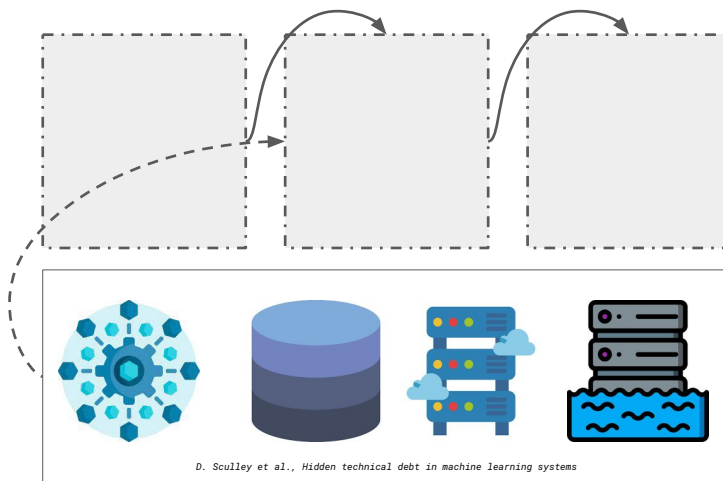
### 2.2 CANCEL shome2022robust

As seen in Figure 4 & 5, there may be several sources of data for the model. For instance, there may be a web service which is continually scrapping the internet for data, or we may have data stored in a database, a data warehouse or data lake. In addition, we may have several other systems with or without ML components which our system communicates with. For instance, our system may rely on data from another service. In return, other services may depend on the predictions from our system.



7

Figure 4: Tangled Systems in ML 1



8

Figure 5: Tangled Systems in ML 2

Thus ML pipelines are inherently complex and tangled. They consist of several stages and the ML model work tends to be a small part of a much larger system. A small change or bug in any of the stages ripples throughout the entire pipeline. Therefore, we cannot make implicit assumptions regarding the quality of the data.

Although data validation libraries exist, currently there are no tools to aid practitioners in writing robust data validation rules. The motivation here is that data will constantly evolve and change as a reflection of the real world. However, this requires a lot of upkeep of data validation rules which are brittle and keep triggering false alerts. Put another way, we want to reduce false positive alerts to on-call engineers to reduce waste of human time & effort.

We started by analysing the data validation rules produced by the Tensorflow Data Validation (TFDV) tool on 218 public datasets. TFDV flagged 74 datasets with anomalies out of 218. We manually analysed the top 5 datasets with the most number of anomalies. However, the results were not fruitful from a research perspective. The premise of TFDV is that data-centric work requires human-in-the-loop approach. The tool provides a basic set of data validation rules, but the practitioners are still expected to write custom validation rules based on domain knowledge.

Next, we tried searching for existing data validations rules from public ML projects. This yet again was unfruitful since there is a lack of public projects which utilise this tool. Projects that do use data validation perhaps cannot make their validation rules public due to privacy & legal obligations.

Our final, albeit unsuccessful attempt was to work with external collaborators in possession of a ML pipeline running in production. We reached out to several researchers & the Dutch Ministry of Transportation & Environment.

### 2.3 TODO shome2022qualitative

Prior work have conducted analysis of data transformations & their impact on the fairness of the ML pipeline [1]. Work also exists on analysing impact of bias mitigation strategies at different stages of the pipeline and their impact on the overall fairness of the pipeline [4].

However, no work has been done to identify how much of the fairness issues identified using both the model & the data, can be identified only using the data. In other words, when can we adopt a data-centric approach to testing for fairness? And when do we require a model as well to make sure our predictions are not biased?

Although our methodology will be similar to the existing work (primarily that of [1]), our research question is entirely different. We are hoping to contribute a systematic methodology for testing for fairness in ML pipelines. With our work, practitioners can choose the best testing strategy given a fairness issue.

## 2.4 LATER shome202xtbd

A natural extension to shome2022qualitative is to conduct a qualitative analysis of bias mitigation techniques. The objective here would be to develop a systematic approach to fixing fairness issues in ML pipelines. Given a fairness issue, we want to aid practitioners in determining the best strategy for fixing the issue.

## 3 Supplemental Contributions

This section presents some contributions I have made outside of my immediate field of research.

- I was a student volunteer at ICSE '22
- I was a teaching assistant for the REMLA course this year. Besides supervising a team of students I also gave a guest lecture on Data Validation for ML.
- I have open sourced the data smells catalogue. I am (slowly) working towards setting up a Github repository to accept contributions.
- I have given a few talks to the group on other topics of interest.
- Caro & I organised SERG's very first Bib meetup where we discussed tools & techniques to manage information, notes & knowledge. Our very own Diomidis was the guest speaker at this event.

## 4 Appendix

### 4.1 DE course overview

Figure 6 presents an overview of the DE courses completed so far. I have completed half of the required credits from the research competency and transferable skill categories. Although I am lagging behind in the discipline

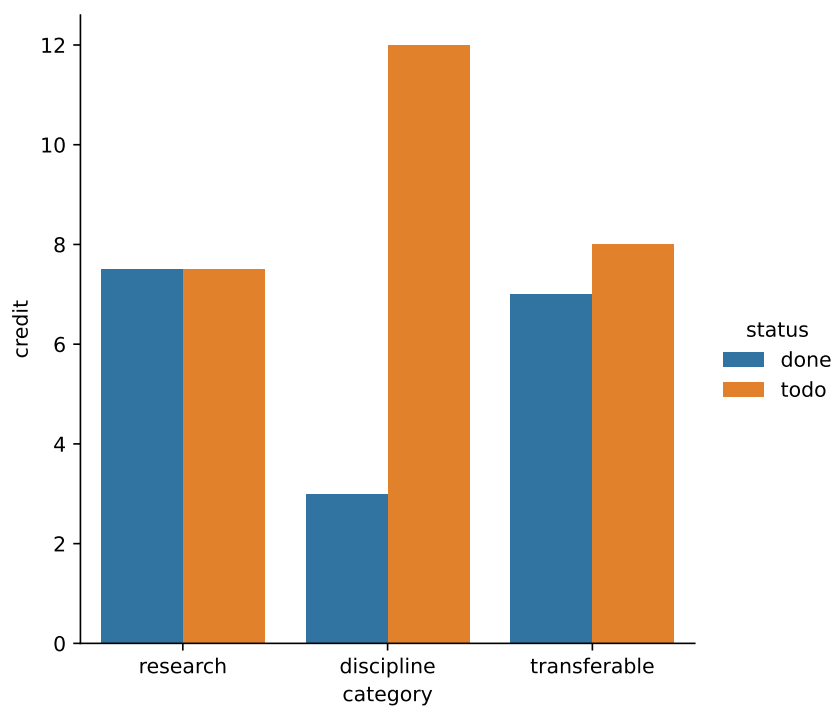


Figure 6: DE Course Overview



skills category, I intend to catch up by attending summer/winter schools in the second half of my Phd.

## 4.2 Data Management Plan

DATASETS	STORAGE	ORGANISATION	DOCUMENTATION	METADATA	FILE FORMAT	ACCESS	PUBLICATION
<b>CODE</b> Code for tools/reviewers are on github. I will be releasing code on github for my research. I have to share with project members, peers & other researchers. My code will be released to be public.	Version controlled using git. I have a local folder named 'code' and I have a remote on Github named 'code'. I will be releasing code on github for my research. I have to share with project members, peers & other researchers. My code will be released to be public.	Follows the folder & file naming convention outlined by the program. I will be releasing code on github for my research. I have to share with project members, peers & other researchers. My code will be released to be public.	Background: Installation instructions, project structure & other general information provided in a README file. I follow the guidelines outlined in the README file. I will be releasing code on github for my research. I have to share with project members, peers & other researchers. My code will be released to be public.	hosted by git & Github.	Code written in Python & documents on text or pdf format. All open formats.	Access to supervisors & project members prior to publication. Open access after publication.	Published on Github under CC BY 4.0 or MIT license.
<b>REPORT</b> I will be releasing reports on github for my research. I have to share with project members, peers & other researchers. My reports will be released to be public.	Version controlled using git. I have a local folder named 'reports' and I have a remote on Github named 'reports'. I will be releasing reports on github for my research. I have to share with project members, peers & other researchers. My reports will be released to be public.	Top level folder named 'reports'. I have a top level folder named 'reports'. I will be releasing reports on github for my research. I have to share with project members, peers & other researchers. My reports will be released to be public.	Project structure & build instructions in a README file. I will be releasing reports on github for my research. I have to share with project members, peers & other researchers. My reports will be released to be public.	hosted by git & Github.	Source written in latex & compiled into pdf with open formats.	Access to supervisors & project members prior to publication. Open access after publication.	Published in conferences & journals. Access depends on publisher.
<b>PUBLIC DATASETS</b> I will be releasing public datasets on github for my research. I have to share with project members, peers & other researchers. My public datasets will be released to be public.	Stored locally or remotely on a public dataset provider. I will be releasing public datasets on github for my research. I have to share with project members, peers & other researchers. My public datasets will be released to be public.	I primarily work with public datasets which are stored in the form of folders named 'raw'. The files are named using the format: 'YYYY-MM-DD-NAME'. I will be releasing public datasets on github for my research. I have to share with project members, peers & other researchers. My public datasets will be released to be public.	Data collection process (e.g. sample source control). I used to document some statistical properties of the datasets in the methods used to analyse data.	None of collection, shape, author, license, version, M, task, device, modified time, etc.	NA	NA	NA
<b>META-DATA</b> I will be releasing meta-data on github for my research. I have to share with project members, peers & other researchers. My meta-data will be released to be public.	Stored locally or remotely on a public dataset provider. I will be releasing meta-data on github for my research. I have to share with project members, peers & other researchers. My meta-data will be released to be public.	Top level folder named 'raw'. For the different versions I use, I use the format: 'YYYY-MM-DD-NAME'. I will be releasing meta-data on github for my research. I have to share with project members, peers & other researchers. My meta-data will be released to be public.	I analyse this in notebooks which document the characteristics, transformations & my results regarding the dataset.	Date of creation, date of last modification, version, shape, etc.	Stored as csv files, open.	Access to supervisors & project members. May also be made public.	Published in <a href="#">Github</a> or <a href="#">Zenodo</a> under CC BY 4.0 license with public DOI.
<b>COMPUTATIONAL NOTEBOOKS</b> I will be releasing computational notebooks on github for my research. I have to share with project members, peers & other researchers. My computational notebooks will be released to be public.	Stored locally or remotely on a public dataset provider. I will be releasing computational notebooks on github for my research. I have to share with project members, peers & other researchers. My computational notebooks will be released to be public.	Top level folder named 'notebooks'. The files are named using the format: 'YYYY-MM-DD-NAME'. I will be releasing computational notebooks on github for my research. I have to share with project members, peers & other researchers. My computational notebooks will be released to be public.	These are self documenting 'raw' files of code is accompanied by my description of the characteristics, transformations & my interpretation of the results.	Author, date of creation, date of last modification, etc.	Stored as ipynb files which are a standard format. I will be releasing computational notebooks on github for my research. I have to share with project members, peers & other researchers. My computational notebooks will be released to be public.	Access to supervisors & project members. May also be made public.	Published in <a href="#">Github</a> or <a href="#">Zenodo</a> under CC BY 4.0 license with public DOI.
<b>FRAGMENTS/INTERVIEWS</b> I will be releasing fragments/interviews on github for my research. I have to share with project members, peers & other researchers. My fragments/interviews will be released to be public.	Stored locally and additional backup on a drive and external SD. I will be releasing fragments/interviews on github for my research. I have to share with project members, peers & other researchers. My fragments/interviews will be released to be public.	No prior experience but a good place to start would be 'raw'. The files are named using the format: 'YYYY-MM-DD-NAME'. I will be releasing fragments/interviews on github for my research. I have to share with project members, peers & other researchers. My fragments/interviews will be released to be public.	File structure & naming convention in a README file. This can also document the interview procedure for future reference.	Name, location, date, time, occupation, country, etc.	Video stored as mp4 files. I will be releasing fragments/interviews on github for my research. I have to share with project members, peers & other researchers. My fragments/interviews will be released to be public.	Access to supervisors & project members. May also be made public.	Stored in a Drive. Restricted access can be provided through the <a href="#">Google Drive</a> sharing settings. I will be releasing fragments/interviews on github for my research. I have to share with project members, peers & other researchers. My fragments/interviews will be released to be public.
<b>PERSONAL NOTES</b> I will be releasing personal notes on github for my research. I have to share with project members, peers & other researchers. My personal notes will be released to be public.	Stored locally & backed up on a private cloud storage provider. I will be releasing personal notes on github for my research. I have to share with project members, peers & other researchers. My personal notes will be released to be public.	I use <a href="#">OneDrive</a> for notes. All notes are stored in a folder named 'notes'. I will be releasing personal notes on github for my research. I have to share with project members, peers & other researchers. My personal notes will be released to be public.	NA	NA	Stored as mp4 files which are a standard format. I will be releasing personal notes on github for my research. I have to share with project members, peers & other researchers. My personal notes will be released to be public.	Closed.	May inspire blog posts. Published on my website under the <a href="#">Creative Commons License</a> .

Figure 7: Data Flow Map

Figure 7 presents the data flow map which provides an overview of how various data generated during the course of this Phd will be handled.

## References

- [1] BISWAS, S., WARDAT, M., AND RAJAN, H. The art and practice of data science pipelines: A comprehensive study of data science pipelines in theory, in-the-small, and in-the-large. In *2022 IEEE/ACM 44th International Conference on Software Engineering (ICSE) (2022)*, IEEE, pp. 2091–2103.
- [2] BOSCH, J., OLSSON, H. H., AND CRNKOVIC, I. Engineering ai systems: A research agenda. In *Artificial Intelligence Paradigms for Smart Cyber-Physical Systems*. IGI Global, 2021, pp. 1–19.
- [3] CRAWFORD, K. *The Atlas of AI*. Yale University Press, 2021.
- [4] GHAI, B., MISHRA, M., AND MUELLER, K. Cascaded debiasing: Studying the cumulative effect of multiple fairness-enhancing interventions. *arXiv preprint arXiv:2202.03734* (2022).

- [5] HUTCHINSON, B., SMART, A., HANNA, A., DENTON, E., GREER, C., KJARTANSSON, O., BARNES, P., AND MITCHELL, M. Towards accountability for machine learning datasets: Practices from software engineering and infrastructure. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (2021)*, pp. 560–575.
- [6] KSHIRSAGAR, M., ROBINSON, C., YANG, S., GHOLAMI, S., KLYUZHIN, I., MUKHERJEE, S., NASIR, M., ORTIZ, A., OVIEDO, F., TANNER, D., ET AL. Becoming good at ai for good. *arXiv preprint arXiv:2104.11757* (2021).
- [7] NG, A. A chat with andrew on mlps: From model-centric to data-centric ai, 2021. Accessed on [2022-01-17 Mon].
- [8] SAMBASIVAN, N., KAPANIA, S., HIGHFILL, H., AKRONG, D., PARITOSH, P., AND AROYO, L. M. “everyone wants to do the model work, not the data work”: Data cascades in high-stakes ai. In *proceedings of the 2021 CHI Conference on Human Factors in Computing Systems (2021)*, pp. 1–15.
- [9] SATO, D., WILDER, A., AND WINDHEUSER, C. Continuous delivery for machine learning.
- [10] SCULLEY, D., HOLT, G., GOLOVIN, D., DAVYDOV, E., PHILLIPS, T., EBNER, D., CHAUDHARY, V., YOUNG, M., CRESPO, J.-F., AND DENNISON, D. Hidden technical debt in machine learning systems. *Advances in neural information processing systems 28* (2015), 2503–2511.
- [11] ZHANG, J. M., HARMAN, M., MA, L., AND LIU, Y. Machine learning testing: Survey, landscapes and horizons. *IEEE Transactions on Software Engineering* (2020).