# TOWARDS REPRODUCIBLE ML RESEARCH IN INFORMATION RETRIEVAL

Ana Lucic, Maurits Bleeker, Maarten de Rijke, Koustuv Sinha, Sami Jullien, Robert Stojnic

July 11, 2022 – 13.30–17.00

**SIGIR 2022** 



#### Yann LeCun @ylecun · Apr 3, 2020

The Transformer-XL results from Google Brain on language modeling could not be reproduced by some top NLP researchers (and the authors are not helping).

@srush\_nlp offers a bounty for whoever can reproduce the results.
(I assume the authors are excluded from the challenge!).

Sasha Rush @srush\_nlp · Apr 2, 2020 Open-Science NLP Bounty: (\$100 + \$100 to charity)

Task: A notebook demonstrating experiments within 30(!) PPL (<84) of this widely cited LM baseline on PTB / WikiText-2 using any non-pretrained, word-only Transformer variant.

Context: twitter.com/Tim\_Dettmers/s...

Show this thread

Model	#Param	PPL
Inan et al. (2016) - Tied Variational LSTM	24M	73.2
Zilly et al. (2016) - Variational RHN	23M	65.4
Zoph and Le (2016) - NAS Cell	25M	64.0
Merity et al. (2017) - AWD-LSTM	24M	58.8
Pham et al. (2018) - Efficient NAS	24M	58.6
Liu et al. (2018) - Differentiable NAS	23M	56.1
Yang et al. (2017) - AWD-LSTM-MoS	22M	55.97
Melis et al. (2018) - Dropout tuning	24M	55.3
Ours - Transformer-XL	24M	54.52

### TUTORIAL OVERVIEW

#### • Part I: Introduction to Reproducibility

• ML reproducibility crisis, examples from non-CS fields, how to conduct reproducible research

#### • Part 2: Reproducibility in IR

• Reproducibility challenges in IR, reproducibility failures in IR, reproducibility tracks

#### • Part 3: Mechanisms for Reproducibility

- Papers with Code, ML Reproducibility Challenge, useful tools and libraries
- Part 4: Reproducibility as a Teaching Tool
  - How to incorporate an ML reproducibility project into a course

### **TEACHING TEAM**



Ana Lucic



Sami Jullien



**Maurits Bleeker** 



Koustuv Sinha



Maarten de Rijke



Robert Stojnic

#### **PREVIOUS TEACHING TEAM**



Ana Lucic University of Amsterdam



Maurits Bleeker

University of Amsterdam



Jesse Dodge





Samarth Bhargav

University of Amsterdam



Sasha Luccioni





Jessica Zosa Forde





**Robert Stojnic** 



Koustuv Sinha



# I – INTRODUCTION TO REPRODUCIBILITY

Ana Lucic

### OVERVIEW

- I. Motivation
- 2. Reproducibility Crisis in ML
- 3. Reproducibility in Non-CS Fields
- 4. Conducting Reproducible Research



"At the very foundation of scientific inquiry is the process of specifying a hypothesis, running an experiment, analyzing the results and drawing conclusions"

"Scientists have used this process to build our collective understanding of the natural world and the laws that govern it. However, for the findings to be valid and reliable, it is important that the experimental process be repeatable, and yield consistent results and conclusions"



solid: submitted; dashed: accepted; red: after 2016

- As a field, we've made considerable progress by increasing the amount of computation used in our experiments:
  - Better performance
  - Easier to explore ideas
- This has also come with some challenges:
  - Running baselines can be very expensive
  - Results are not always reproducible

How can we mitigate the challenges of bigger, more complex models without reducing the benefits?

How can we mitigate the challenges of bigger, more complex models without reducing the benefits?

#### In this tutorial, we focus on the challenge of <u>ensuring</u> <u>research results are reproducible</u>

### TUTORIAL OVERVIEW

- I. Introduction to Reproducibility
- 2. Reproducibility in IR
- 3. Mechanisms for Reproducibility
- 4. Reproducibility as a Teaching Tool



Baker. 2016. Is there a reproducibility crisis? Nature.











# I.b – REPRODUCIBILITY IN ML

## ACM DEFINITIONS (v1.1)

- **Repeatable:** a researcher can obtain the same results for their own experiment under exactly the same conditions, i.e., they can reliably repeat their own experiment ("Same team, same experimental setup")
- **Reproducibility:** a different researcher can obtain the same results for an experiment under exactly the same conditions and using exactly the same artifacts, i.e., another independent researcher can reliably repeat an experiment of someone other than herself ("Different team, same experimental setup") [this was called *replicability* in v1.0 of ACM definitions]
- Replicability: a different researcher can obtain the same results for an experiment under different conditions and using their self-developed artifacts ("Different team, different experimental setup") [this was called *reproduciblity* in v1.0 of ACM definitions]

### **NEURIPS DEFINITIONS**

- **Reproducible:** same conclusions are drawn when re-doing an experiment with the same data and same analytical tools
- **Replicable:** same conclusions are drawn when re-doing an experiment with a different dataset, but the same tools
- **Robust:** same conclusions are drawn when re-doing an experiment with the same data but different tools (i.e., different code implementations)
- **Generalizable:** same conclusions are drawn when re-doing an experiment with different data and different tools.

#### **NEURIPS DEFINITIONS**



#### **REPRODUCIBILITY CRISIS IN ML**

# **Code break**

In a survey of 400 artificial intelligence papers presented at major conferences, just 6% included code for the papers' algorithms. Some 30% included test data, whereas 54% included pseudocode, a limited summary of an algorithm.



#### **REPRODUCIBILITY CRISIS IN ML**

#### Code and Data Associated with this Article

orXiv Links to Code & Data (What is Links to Code & Data?)

#### **Official Code**

No official code found; you can submit it here

#### **Community Code**

[111] 5 code implementations (in PyTorch and TensorFlow)

#### **Datasets Used**

**OpenAl Gym** 853 papers also use this dataset

MuJoCo 831 papers also use this dataset

- Since 2018, we've made some progress
- Many conferences strongly encourage or even require code submissions
- Can get links to code repositories and datasets through arXiv thanks to Papers with Code
- Reproducibility checklists at conferences

### COMMON REPRODUCIBILITY ISSUES IN ML

- Lack of access to the same training data, differences in data distribution
- Misspecification or under-specification of the model or training procedure
- Lack of availability of the code necessary to run the experiments, or errors in the code
- Under-specification of the metrics used to report results
- Improper use of statistics to analyze results
- Selective reporting or over-claiming of results

## QUESTIONS?

# I.c – REPRODUCIBILITY IN NON-CS FIELDS

## PSYCHOLOGY

- The Open Science Collaboration conducted 100 replications of studies from 3 psychology journals
  - In total, there are 270 authors on the paper published in Science
- Found a significant proportion of replications produced weaker evidence despite using materials provided by authors
- Mean effect size of replication was found to be half of the original
  - Original: 97% significant (p< 0.05) vs Study: 36%

### **BIOMEDICAL SCIENCES**

- Clinical trials in oncology have some of the highest failure rates in comparison to other therapeutic areas
- Begley and Lee (2012) claim this is due to the lack of robustness in preclinical trials i.e., drug development
- Out of 53 "landmark" studies, only 6 could be reproduced
- Non-reproducible papers are still heavily cited since they are considered to be "part of the literature", contributing to failing clinical trials

#### **BIOMEDICAL SCIENCES**

#### **REPRODUCIBILITY OF RESEARCH FINDINGS**

Preclinical research generates many secondary publications, even when results cannot be reproduced.

Journal impact factor	Number of articles	Mean number of citations of non-reproduced articles*	Mean number of citations of reproduced articles
>20	21	248 (range 3–800)	231 (range 82–519)
5–19	32	169 (range 6–1,909)	13 (range 3–24)

Results from ten-year retrospective analysis of experiments performed prospectively. The term 'non-reproduced' was assigned on the basis of findings not being sufficiently robust to drive a drug-development programme. \*Source of citations: Google Scholar, May 2011.

#### **BIOMEDICAL SCIENCES**

Recommendations proposed by Begley and Lee (2012):

- Require reporting on negative findings
- Encourage reporting on alternative findings that contradict existing work
- Implement transparent mechanisms for reporting unethical practices
- Increase dialogue between physicians, scientists, patient advocates and patients
- Recognize high-quality teaching and mentoring as valuable
- Funding organizations should facilitate development and access to new tools

## CANCER BIOLOGY

Errington et al (2020) conduct a reproduction of 193 experiments from 53 high impact papers in preclinical cancer biology:

- Only 50/193 experiments from 23 papers were reproduced
- Data was publicly accessible for 4 of 193 papers
- Authors would not share data for 68% of papers
- 32% of authors were rated as "not at all helpful" by researchers reproducing their experiments
- 67% of protocols described in papers needed modifications
  - Only 41% of those modifications could be implemented

#### CANCER BIOLOGY



### ECONOMICS

- Camerer et al (2016) analyze 18 studies in economics:
- They find that 61% of studies detect the original effect size in the same direction at alpha = 0.05
- However, the replicated effect size is 66% of the original, on average


# I.d – CONDUCTING REPRODUCIBLE RESEARCH

# CONDUCTING REPRODUCIBLE RESEARCH

- 1. Hypothesis testing
- 2. Randomness
- 3. Statistical testing
- 4. Open-source code
- 5. Model cards
- 6. Datasheets

# HYPOTHESIS TESTING

- In ML/IR, we often get started with running experiments right away due to the low barrier to entry, which can result in:
  - Unclear research questions
  - Unclear conclusions
  - Wasted time, effort and computation power
- Formulating (some version of) the RQs before starting with experimentation can help alleviate some of these issues

# RANDOMNESS

Deep Neural Networks display highly non-convex loss surfaces and therefore the performance of a model depends on several factors:

- Specific hyperparameters
- Dropout applied during training
- Weight initialization
- Order of the training data
- Randomly sampled data augmentations

It is important identify all sources of potential randomness in order to try to compensate for them in your experiments

# STATISTICAL TESTING

- Comparing the means of two models is not enough to conclude model A is better than B
- It is important to choose the appropriate statistical test to determine whether or not your results are significant. Some resources:
  - Ulmer et al. 2022. Deep-Significance: Easy and Meaningful Statistical Significance Testing in the Age of Neural Networks.
  - Dror et al. 2019. Deep Dominance: How to Properly Compare Deep Neural Models.

### STATISTICAL TESTING

- Scenario I: Comparing multiple runs of two models
  - Scores from a model **A** and a baseline **B** on a dataset, stemming from N model runs with different random seeds
  - Comparing multiple runs will *always* be preferable
- Scenario 2: Comparing multiple runs across datasets
  - When comparing models across datasets, formulate one null hypothesis per dataset
  - N model runs with different random seeds

# STATISTICAL TESTING

- Scenario 3: Comparing sample-level scores
  - If only one run is available, comparing sample-wise score distribution can be an option
- Scenario 4: Comparing more than two models
  - For instance, for three models, we can create a matrix 3x3
- The framework by Ulmer et al. 2022 makes use of the Almost Stochastic Order (ASO) test
  - Expresses the amount of violation of stochastic order

## **OPEN-SOURCE CODE**

- When possible, it is beneficial to open source your code and data in order to promote open and reproducible science
- Templates such as the ML Code Completeness Checklist can help you arrange your repository before publishing it publicly
  - More details in Part 3 of the tutorial
- Open-source code provides insights into:
  - The underlying implementation of a formal idea
  - Many hyperparameters and minor details that are not discussed in the paper

#### MODEL CARDS

#### Model Card - Toxicity in Text

#### **Model Details**

- The TOXICITY classifier provided by Perspective API [32], trained to predict the likelihood that a comment will be perceived as toxic.
- Convolutional Neural Network.
- Developed by Jigsaw in 2017.

#### Intended Use

- Intended to be used for a wide range of use cases such as supporting human moderation and providing feedback to comment authors.
- Not intended for fully automated moderation.
- Not intended to make judgments about specific individuals.

#### Factors

• Identity terms referencing frequently attacked groups, focusing on sexual orientation, gender identity, and race.

#### Metrics

 Pinned AUC, as presented in [11], which measures threshold-agnostic separability of toxic and non-toxic comments for each group, within the context of a background distribution of other groups.

#### **Ethical Considerations**

 Following [31], the Perspective API uses a set of values to guide their work. These values are Community, Transparency, Inclusivity, Privacy, and Topic-neutrality. Because of privacy considerations, the model does not take into account user history when making judgments about toxicity.

#### **Training Data**

- Proprietary from Perspective API. Following details in [11] and [32], this includes comments from a online forums such as Wikipedia and New York Times, with crowdsourced labels of whether the comment is "toxic".
- "Toxic" is defined as "a rude, disrespectful, or unreasonable comment that is likely to make you leave a discussion."

#### **Evaluation Data**

- A synthetic test set generated using a template-based approach, as suggested in [11], where identity terms are swapped into a variety of template sentences.
- Synthetic data is valuable here because [11] shows that real data often has disproportionate amounts of toxicity directed at specific groups. Synthetic data ensures that we evaluate on data that represents both toxic and non-toxic statements referencing a variety of groups.

#### **Caveats and Recommendations**

 Synthetic test data covers only a small set of very specific comments. While these are designed to be representative of common use cases and concerns, it is not comprehensive.

## DATASHEETS

Datasheets were proposed as a mechanism to standardize documentation practices for ML datasets. They include ~50 questions on the following topics:

- Motivation
- Composition
- Collection Process
- Preprocessing/cleaning/labelling
- Uses
- Distribution
- Maintenance

#### DATASHEETS

#### Motivation

For what purpose was the dataset created? Was there a specific task in mind? Was there a specific gap that needed to be filled? Please provide a description.

The dataset was created to enable research on predicting sentiment polarity—i.e., given a piece of English text, predict whether it has a positive or negative affect—or stance—toward its topic. The dataset was created intentionally with that task in mind, focusing on movie reviews as a place where affect/sentiment is frequently expressed.<sup>1</sup>

#### Who created the dataset (e.g., which team, research group) and on behalf of which entity (e.g., company, institution, organization)?

The dataset was created by Bo Pang and Lillian Lee at Cornell University.

#### Who funded the creation of the dataset? If there is an associated grant, please provide the name of the grantor and the grant name and number.

Funding was provided from five distinct sources: the National Science Foundation, the Department of the Interior, the National Business Center, Cornell University, and the Sloan Foundation.

#### Any other comments?

None.

#### Distribution

Will the dataset be distributed to third parties outside of the entity (e.g., company, institution, organization) on behalf of which the dataset was created? If so, please provide a description.

Yes, the dataset is publicly available on the internet.

# How will the dataset will be distributed (e.g., tarball on website, API, GitHub)? Does the dataset have a digital object identifier (DOI)? The dataset is distributed on Bo Pang's webpage at Cornell: http://www.cs.cornell.edu/people/pabo/movie-review-data. The dataset does

not have a DOI and there is no redundant archive. When will the dataset be distributed?

The dataset was first released in 2002.

Will the dataset be distributed under a copyright or other intellectual property (IP) license, and/or under applicable terms of use (ToU)? If so, please describe this license and/or ToU, and provide a link or other

access point to, or otherwise reproduce, any relevant licensing terms or ToU, as well as any fees associated with these restrictions.

The crawled data copyright belongs to the authors of the reviews unless otherwise stated. There is no license, but there is a request to cite the corresponding paper if the dataset is used: *Thumbs up? Sentiment classification using machine learning techniques.* Bo Pang, Lillian Lee, and Shivakumar Vaithyanathan. Proceedings of EMNLP, 2002.

#### RECOMMENDATIONS FOR CONDUCTING REPRODUCIBLE RESEARCH

- 1. Formulate hypothesis prior to starting experiments
- 2. Identify appropriate statistical tests
- 3. Identify stochastic components of experiments and account for randomness
- 4. Open-source your code with clear instructions on how to run it
- 5. Clearly document your contribution with a model card and/or a datasheet

# QUESTIONS?

### MOTIVATION

How can we mitigate the challenges of bigger, more complex models without reducing the benefits?

### In this tutorial, we focus on the challenge of <u>ensuring</u> <u>research results are reproducible</u>

# TUTORIAL OVERVIEW

- I. Introduction to Reproducibility
- 2. Reproducibility in IR
- 3. Mechanisms for Reproducibility
- 4. Reproducibility as a Teaching Tool

# II – REPRODUCIBILITY IN IR

Sami Jullien, Maarten de Rijke

## OVERVIEW

- L. Reproducibility challenges in IR
- 2. Reproducibility failures in IR
- 3. Reproducibility tracks and initiatives in IR

# II.a REPRODUCIBILITY CHALLENGES IN IR

# **TEST COLLECTIONS**

#### • Cranfield tradition

- Create "laboratory type situation where, freed as far as possible from the contamination of operational variables, the performance of index languages could be considered in isolation"
- Same set of documents, same set of information needs to be used for each index language, and for the use of both precision and recall to evaluate the effectiveness of the search.
  - Relevance was based on topical similarity where the judgments were made by domain experts
  - Relevance = topic similarity, single set of judgments representative, set of relevant docs is complete
- Almost all sources of variability removed in design of test collections: users, tasks, leaving topics ("statements of information needs") as main source of variability in collection
- Judgments indicating which documents are relevant to which topics

## CRANFIELD SCALING UP: TREC

Launched in 1992:

- Provides document set and a set of topics to the participants.
- Each participant runs topics against documents using their retrieval system, and returns to NIST a ranked list of the top *N* documents per topic
- TREC forms **pools** from participants' submissions, which are judged by relevance assessors
- Submissions are evaluated using resulting relevance judgments; evaluation results returned to participants

# CRANFIELD SCALING UP: TREC

Thirty years later:

- Pioneered use of "pooling" for building large collections
- So far built > 150 test collections for dozens of search tasks
- Hundreds of participant teams world-wide
- Premier venue for determining research methodology
- Model for other efforts in IR and related fields

# TREC

Large number of tasks, partly inspired by available funding, partly by community initiatives

- Different document genres
- Different domains
- Different types of information need
- Different relevance criteria
- Different languages



CI FF

Conference and Labs of the Evaluation Forum

- Running since 2000
- "Promote research, innovation, and development of information access systems with an emphasis on multilingual and multimodal information with various levels of structure"
- **Conference part** with focus on experimentation in IR
- Labs part similar to TREC with very diverse set of tasks
  - CLEF 2022 (Bologna, September 2022) features 14 labs
    - (1) Answer Retrieval for Questions on Math, (2) Large-scale biomedical semantic indexing and question answering, (3) Fighting the COVID-19 Infodemic and Fake News Detection, (4) Cheminformatics, (5) Early risk prediction on the Internet, (6) Named Entity Recognition and Linking in Multilingual Historical Documents, (7) Intelligent Disease Progression Prediction, (8) ImageCLEF: Multimedia Retrieval Challenge, (9) Automatic Wordplay and Humour Translation, (10) Learning to Quantify, (11) Biodiversity identification and prediction challenges, (12) Digital Text Forensics and Stylometry, (13) SimpleText: Automatic Simplification of Scientific Texts, (14) Argument Retrieval

### NTCIR

NTCIR (NII Test Collection for IR Systems) Project

- Running since 1997
- "Evaluation efforts designed to enhance research on diverse information access technologies, including, but not limited to, cross-language and multimedia information access, question-answering, text mining and summarisation, with an emphasis on East Asian languages such as Chinese, Korean, and Japanese, as well as English"
- Runs on an 18 month cycle, 2022 edition just took place (NTCIR-16, June 2022)
- Tasks at NTCIR-16
  - Main: (1) Data Search, (2) Dialogue Evaluation, (3) Investor's and Manager's Fine-grained Claim Detection, (4) Lifelog Access and Retrieval, (5) QA Lab for Politics Information (6) We Want Web 4 with CENTER
  - **Pilot**: (1) Reading Comprehension for Information Retrieval, (2) Real document-based Medical Natural Language Processing, (3) Session Search, (4) Unbiased Learning to Ranking Evaluation Task

#### FIRE

Forum for Information Retrieval Evaluation

- Running since 2008
- "Encourage research in Indian language Information Access technologies by providing reusable large-scale test collections for Indian language IR experiments; Provide a common evaluation infrastructure for comparing the performance of different IR system; Investigate evaluation methods for Information Access techniques and methods for constructing a reusable large-scale data set for Indian language IR experiments"
- Mixture of conference and collaborative benchmarking
- FIRE 2022 (Kolkata, December 2022) features 8 tracks
  - (1) Anaphora Resolution from Social Media Text in Indian Languages, (2) Emotions & Threat Detection in Urdu, (3) Hate Speech and Offensive Content Identification in English and Indo-Aryan Languages, (4) Indian Language Summarization, (5) Information Retrieval from Microblogs during Disasters, (6) Information Retrieval in Software Engineering, (7) Self-reported Mental Disorder Diagnosis, (8) Sentiment Analysis and Homophobia detection of YouTube comments in Code-Mixed Dravidian Languages

### ALIGNING: CENTRE

- CLEF, NTCIR, TREC REproducibility
- Aims:
  - Reproduce best results of best/most interesting systems in previous editions of CLEF/NTCIR/TREC by using standard open source IR systems
  - Contribute back to the community the additional components and resources developed to reproduce the results in order to improve existing open source systems
- Running since TREC 2018

### ADVANTAGES OF COMMUNITY-BASED BENCHMARKING

#### • Improve the **state of the art**

- Make research and research results comparable
- Solidify a research **community** 
  - Create resources together, iterate over successes and failures
- Amortize the costs of infrastructure
  - Collaborative development and analysis
- Facilitate technology transfer
  - Compare using an independent, open yardstick before trying out "at home"
- Establish research methodology
  - Scrutinize it ...

# QUESTIONS ABOUT METHODOLOGY: VALIDITY

#### • Internal validity

- Does way in which study was done allow trustworthy answers to its research questions?
- Examines, e.g., the extent to which systematic error (bias) is present.

#### • External validity

- Can the findings of a study can be generalized to **other contexts**?
  - E.g., other queries, other document collections, other TREC collections, ...

#### • Ecological validity

- Can the results of a study can be generalized to real-life settings?
  - Differs from external validity
  - E.g., from TREC collection to AB test in a production e-commerce setting

#### Where does reproducibility/replicability come in?

- Repeatability: same team, same setup
- Reproducibility: different team, same setup
- Replicability: different team, different setup

### INTERNAL VALIDITY

Dimensions to consider when considering internal validity

- Improper randomization so that you're not really looking at a random sample of users or queries
- Controlled pooling produces unbiased judgment set sufficient for comparative evaluation
- Not running sufficiently many • iterations in an experiment may lead to wrong conclusions
  - After 30,000 impressions no noticeable performance difference between linear and neural online ITR
  - After 1,000,000 there is a difference





#### EXTERNAL VALIDITY

#### • Variability

- How spread out or closely clustered a set of data instances is (queries, users, tasks, documents, ...)
- Bailey et al. consider **variability across users** and especially across different individual query formulations and expectations of quantities of relevant information needing to be found
  - Does the existence of individual variation in initial query formulation for a single information need alter the evaluation of system performance?
  - Is there significant variation among users of the anticipated effort in terms of the number of documents viewed and queries to be issued, and is there a relationship between a user's anticipated effort and the information task complexity?
- Findings
  - Both questions answered affirmatively.
  - The use of multiple queries per topic arising from different searchers provides a more representative characterization of the mapping from information need than just one

### ECOLOGICAL VALIDITY

Do conclusions reached from Cranfield experiments transfer to operational settings?

- Unable to verify conclusions from a laboratory experiment in user studies
- User studies did not show that conclusions from the laboratory test were wrong, simply that the user studies could not detect any differences
- Relevance judges typically do not assess queries and documents that reflect their own information needs, and have to make assumptions about relevance from an assumed users point of view
- Because the true information need can be difficult to assess, this can cause substantial biases

### ECOLOGICAL VALIDITY

To address gap between offline evaluation and true use of IR systems, **online evaluation** has been used to directly measure observable user behavior on alternative systems

- Challenge for online evaluation is to identify metrics that accurately reflect user satisfaction
- CTR, ranks of clicked documents, skips, time-between-revisits, SAT clicks, ...
- Which are indicative of outcomes of AB tests?

# II.b – REPRODUCIBILITY FAILURES IN IR

#### IMPROVEMENTS DON'T ADD UP

Armstrong et al.

- "There is [...] no evidence that ad-hoc retrieval technology has improved during the past decade or more."
- Finding was arrived at by comprehensive longitudinal survey of research papers between 1998 and 2008 from major IR research venues that report results on a diverse range of TREC test collections.
- Analysis points to "selection of weak baselines that can create an illusion of incremental improvement" and "insufficient comparison with previous results"



Internal validity

# IMPROVEMENTS DON'T ADD UP

Armstrong et al.

- How confident are we that a technique that yields an improvement over a weak baseline would also give an improvement over a strong one, and therefore be a worthwhile addition to state of the art systems?
- Recommendations going forward
  - Avoid perverse selection bias where statistically significant improvements can only be obtained in comparisons against weak baselines
  - Adopt a practice of regular longitudinal comparison to ensure measurable progress, or at least prevent the lack of it from going unnoticed → maintain a leaderboard?

validity

### WEAK BASELINES

#### Lin 2018

- A SOTA claim is an informal prerequisite to get a publication in a top IR conference
- Comparing against the same weak baselines results in a lack of a de facto leaderboard
- A frequent behavior in IR is to compare against cherry-picked points of comparison
- "Pick the best implementation [...] **implementations matter, more so than models**, and thus it makes sense to pick the best one"
- Tuning baselines diminishes how great a proposed method is

validity
### ARE WE REALLY MAKING PROGRESS?

A repeat of Armstrong et al., but now for top-*n* recommendation

Do deep learning methods for top-*n* recommendation really outperform simpler methods?

- Systematic analysis of algorithmic proposals for top-*n* recommendation tasks.
- 18 algorithms presented at top-level research conferences in the last years
- Only 7 could be reproduced with reasonable effort
- 6 of those 7 can often be outperformed with simple heuristic methods (nearest-neighbor or graph-based techniques)

Table 1: Reproducible works on deep learning algorithmsfor top-n recommendation per conference series from 2015to 2018.

Conference	Rep. ratio	Reproducible				
KDD	3/4 (75%)	[17], [23], [48]				
RecSys	1/7 (14%)	[53]				
SIGIR	1/3 (30%)	[10]				
WWW	2/4 (50%)	[14], [24]				
Total	7/18 (39%)					
Non-reprodu	Non-reproducible: KDD: [43], RecSys: [41], [6], [38],					
[44], [21], [45], SIGIR: [32], [7], WWW: [42], [11]						

Table 1. Statistics of relevant and reproducible works on deep learning algorithms for *top-n* recommendation per conference series from 2015 to 2018.

Conference	Rep. Setup ratio	Reproducible Setup	Non-Reproducible Setup
KDD	3/4 (75%)	[32], [37], [73]	[67]
IJCAI	5/7 (71%)	[29], [80], [79], [13], [77]	[51], [76]
WWW	2/4 (50%)	[30], [38]	[66], [22]
SIGIR	1/3 (30%)	[21]	[48], [12]
RecSys	1/7 (14%)	[81]	[65], [8], [60], [68], [34], [71]
WSDM	0/1 (0%)		[75]
Total	12/26 (46%)		

Interna validity

# WHY WE SEE PERFORMANCE DIFFERENCES

#### Li et al. 2022

- Even if we do see meaningful performance differences, **make sure we get the analysis right and we give credit where credit is due** and understand why we are doing a better job
- Next basket recommendation (in grocery shopping)
  - Complex NBR models were in some **but not all** cases able to beat very simple baselines (top frequency, personalized top frequency)
    - Weak or missing baselines, the use of different datasets in different papers, and of non-standard metrics
  - But a close look at the problem space shows that the real algorithmic challenge is one of repetition as most people consume roughly the same, over and over again (in grocery)
  - "when a measure becomes a target, it ceases to be a good measure"
  - So this is not about similarity and complex representation learning, but about understanding temporal patterns

validity

# II.c – REPRODUCIBILITY TRACKS AND INITIATIVES IN IR

# EXISTING REPRODUCIBILITY TRACKS IN IR

- Pressure to publish is often seen as a driver towards non-reproducible research
- Thus, we need reproducibility tracks to evaluate if the progress we are making is applicable for different tasks
- Those tracks are still a minor part of the conferences
- There are several existing reproducibility tracks at IR conferences, most of which focus on replicability + reproducibility as defined by ACM
  - e.g., ECIR, RecSys, SIGIR, Sim4IR workshop

# EXISTING REPRODUCIBILITY TRACKS IN IR

Existing reproducibility tracks at IR conferences:

- ECIR:
  - Track introduced in 2015
  - 2022 version had 11 accepted papers
- RecSys:
  - Track introduced in 2020
  - 2021 version had 3 accepted papers (2022 not announced yet)
- SIGIR:
  - Track introduced in 2022 (!)
  - Strong emphasis on generalizability of lessons learned
  - 7 accepted papers

### SIM4IR

Balog et al. 2021

- Different scenarios can be explored to determine the effects of the simulations' parameters. These scenarios can be run in such a way to ensure reproducible results, with all this being achieved at a low cost to researchers.
- "It is important to note that simulators do not need to be perfect mirrors of human behaviour, but instead simply need to be 'good enough'."
- **"The main requirement is reproducibility".** Non-deterministic simulators should come with random seed numbers to ensure repeatable trajectories.

validity

# **REPRODUCIBILITY FRAMEWORKS**

- MultiReQA: Cross-domain evaluation for retrieval question answering models
  - https://github.com/google-research-datasets/MultiReQA
- KILT: Benchmark for knowledge intensive language tasks
  - https://github.com/facebookresearch/KILT
- BEIR: Heterogeneous benchmark containing diverse IR tasks
  - https://github.com/beir-cellar/beir



Figure 1: An overview of the diverse tasks and datasets in BEIR benchmark.

External

## MOTIVATION

How can we mitigate the challenges of bigger, more complex models without reducing the benefits?

## In this tutorial, we focus on the challenge of <u>ensuring</u> <u>research results are reproducible</u>

# QUESTIONS?

# TUTORIAL OVERVIEW

- I. Introduction to Reproducibility
- 2. Reproducibility in IR
- 3. Mechanisms for Reproducibility
- 4. Reproducibility as a Teaching Tool

# III – MECHANISMS FOR REPRODUCIBILITY

Koustuv Sinha, Robert Stojnic

## OVERVIEW

- I. Papers with Code
- 2. Reproducibility Challenge
- 3. Reproducibility Checklists
- 4. Useful Tools and libraries



• Goal: Track all artefacts in ML, create positive incentives for sharing

[00]	Search Q	Browse State-of-the-Art	Datasets	Methods	More $\vee$	We are hiring!	У.	1	Sign In
	Top & Social 🔅 New 🖓	' Greatest							
	Trending Research						Subscribe		
		MVSTER: Epipolar	Transforme	er for Efficie	ent Multi-V	/iew Stereo	★ 52		
		C Jeffwang987/mvster •	STER, which lev	verages the proj	oosed epipolar	Transformer to learn both 2	1.29 stars / hour		
		semantics and 3D spatial a	associations effic	ciently.			Paper		



• Largest database of papers curated with their code

Code		🖉 Edit
carolineec/EverybodyDanceNow          O official	★ 508	O PyTorch
O Lotayou/everybody_dance_now_pytorch	* 256	O PyTorch
♥ VisiumCH/AMLD2020-Dirty-Gancing ➡ Quickstart in <sup>∞</sup> Colab	★ 17	O PyTorch
O wjy5446/pytorch-everybody-dance-now	★ 9	O PyTorch
O Novemser/deep-imitation	★ 9	O PyTorch
See all 14 implementations		

#### Largest database of datasets, tracking their usage

#### ImageNet

#### Introduced by Jia Deng et al. in ImageNet: A large-scale hierarchical image database

The **ImageNet** dataset contains 14,197,122 annotated images according to the WordNet hierarchy. Since 2010 the dataset is used in the ImageNet Large Scale Visual Recognition Challenge (ILSVRC), a benchmark in image classification and object detection. The publicly released dataset contains a set of manually annotated training images. A set of test images is also released, with the manual annotations withheld. ILSVRC annotations fall into one of two categories: (1) image-level annotation of a binary label for the presence or absence of an object class in the image, e.g., "there are cars in this image" but "there are no tigers," and (2) object-level annotation of a tight bounding box and class label around an object instance in the image, e.g., "there is a screwdriver centered at position (20,25) with width of 50 pixels and height of 30 pixels". The ImageNet project does not own the copyright of the images, therefore only thumbnails and URLs of images are provided.

- Total number of non-empty WordNet synsets: 21841
- Total number of images: 14197122

Homepage

- Number of images with bounding box annotations: 1,034,908
- Number of synsets with SIFT features: 1000
- Number of images with SIFT features: 1.2 million

Source: 🗅 ImageNet Large Scale Visual Recognition Challenge



Source: https://cs.stanford.edu/people/ka





🖻 Edit

87



• Largest database of results from published papers

#### Image Classification on ImageNet





Integrated with:

- arXiv
- ACL anthology
- OpenReview

Bibliographic Tools	Code & Data	Demos	Related Papers	About arXivLabs			
Code and Data Associated with this Article							
arXiv Links to Code & Data (What is Links to Code & Data?)							
Official Code							
https://github.c	https://github.com/carolineec/EverybodyDanceNow						
Community Cod	e						
[IIII] 13 code implen	[IIII] 13 code implementations (in PyTorch)						
Datasets Used							
Everybody Da     ★ introduced in th     7 papers also use t	nce Now is paper :his dataset						



Reproducibility reports shown next to original papers

#### Deep Fair Clustering for Visual Learning

CVPR 2020 · Peizhao Li, Han Zhao, Hongfu Liu · 🖻 Edit social preview

Fair clustering aims to hide sensitive attributes during data partition by balancing the distribution of protected subgroups in each cluster. Existing work attempts to address this problem by reducing it to a classical balanced clustering with a constraint on the proportion of protected subgroups of the input space...

#### 🖾 PDF 📑 Abstract

#### **Reproducibility Reports**

#### Jan 31 2021

#### [Re] Deep Fair Clustering for Visual Learning

RC 2020 · Pauline Baanders, Chris Al Gerges, Nienke Reints, Tobias Teule

For the MNIST-USPS dataset, we report similar accuracy and NMI values that are within 1.2% and 0.5% of the values reported in the original paper. However, the balance and entropy differed significantly, where our results were within 73.1% and 30.3% of the original values respectively. For the Color Reverse MNIST dataset, we report similar values on accuracy, balance and entropy, which are within 5.3%, 2.6% and 0.2% respectively. Only the value of the NMI differed significantly, name within 12.9% of the original value In general, our results still support the main claim of the original paper, even though on some metrics the results differ significantly.



Collated resources for publishing research code

aperswith	ncode / <mark>releasing-researc</mark> l	n-code (Public)	ি Constant Section Se	53 - 😵 Fork 572	🛉 Starred 1.9k 👻
<> Code ⊙	) Issues 2 🕅 Pull requests	🕑 Actions 🗄 Projects 🖽 Wiki	i 🕕 Security 🗠 Insights 🕸 Settin	gs	
٤ <sup>9</sup> maste	ipinic Update README.md ebooks	Fix graph Update README.md	Go to file     Add file •     Code •       a5b2c85 on Mar 19, 2021     ③ 120 commits       2 years ago       2 years ago       2 years ago	About Tips for releasing researc Machine Learning (with or 2020 recommendations) machine-learning awesom peuring peuring-2020	र्छ h code in fficial NeurIPS ne-list
Lice REAL	ENSE NDME.md	Create LICENSE Update README.md	2 years ago 14 months ago	따 Readme 한 MIT License ☆ 1.9k stars	
<b>Tip</b>	DIated best practices from mo	Research Code	row official guidelines at NeurIPS	So so watching So so	



ML Code Completeness Checklist (Robert Stojnic, 2020)



- 1. **Dependencies** does a repository have information on dependencies or instructions on how to set up the environment?
- 2. Training scripts does a repository contain a way to train/fit the model(s) described in the paper?
- 3. Evaluation scripts does a repository contain a script to calculate the performance of the trained model(s) or run experiments on models?
- 4. **Pretrained models** does a repository provide free access to pretrained model weights?
- 5. **Results** does a repository contain a table/plot of main results and a script to reproduce those results?

# QUESTIONS?

### **REPRODUCIBILITY CHECKLISTS**

- ML Reproducibility Checklist (Joelle Pineau, 2018)
- Minimal information that should be in a manuscript
- Not necessarily exhaustive
- Part of guidelines for major conferences (NeurIPS, ICML, ICLR)

#### The Machine Learning Reproducibility Checklist (v2.0, Apr.7 2020)

For all models and algorithms presented, check if you include:

- A clear description of the mathematical setting, algorithm, and/or model.
- A clear explanation of any assumptions.
- An analysis of the complexity (time, space, sample size) of any algorithm.
- For any theoretical claim, check if you include:
- A clear statement of the claim.
- A complete proof of the claim.
- For all datasets used, check if you include:
- The relevant statistics, such as number of examples.
- The details of train / validation / test splits
- An explanation of any data that were excluded, and all pre-processing step.
- A link to a downloadable version of the dataset or simulation environment.
- For new data collected, a complete description of the data collection process, such as instructions to annotators and methods for quality control.
- For all shared code related to this work, check if you include:
- Specification of dependencies.
- Training code.
- Evaluation code.
   Pre-trained model(s)
- Pre-trained model(s).
- □ README file includes table of results accompanied by precise command to run to produce those results.
- For all reported experimental results, check if you include:
- The range of hyper-parameters considered, method to select the best hyper-parameter configuration, and specification of all hyper-parameters used to generate results.
- The exact number of training and evaluation runs.
- A clear definition of the specific measure or statistics used to report results.
- A description of results with central tendency (e.g. mean) & variation (e.g. error bars).
- The average runtime for each result, or estimated energy cost.
- A description of the computing infrastructure used.

Reproduced from: www.cs.mcgill.ca/~jpineau/ReproducibilityChecklist-v2.0.pdf

- Started 2018, till date five editions: ICLR 2018, ICLR 2019, NeurIPS 2019, RC 2020, RC 2021
- Task: Choose a submitted paper from a conference, reproduce the central claim of the paper

#### ML Reproducibility Challenge 2021 Edition

for papers published in:





#### **Best Paper Award**

▶ Reproducibility Study of "Counterfactual Generative Networks", Piyush Bagad, Jesse Maas, Paul Hilders, Danilo de Goede, Forum, Original Paper (ICML 2021)

#### Outstanding Paper Awards

▶ [Re] Learning to count everything, Matija Teršek, Domen Vreš, Maša Kljun, Forum, Original Paper (CVPR 2021)

▶ [RE] An Implementation of Fair Robust Learning , Ian Hardy, Forum, Original Paper (ICML 2021)

► Strategic classification made practical: reproduction, *Guilly Kolkman, Maks kulicki, Jan Athmer, Alex Labro*, Forum, Original Paper (ICML 2021)

► On the reproducibility of "Exacerbating Algorithmic Bias through Fairness Attacks", Andrea Lombardo, Matteo Tafuro, Tin Hadži Veljković, Lasse Becker-Czarnetzki, Forum, Original Paper (AAAI 2021)



Reproducibility Reports accepted to MLRC 2021 by conference

#### Volume 7 (2021)

#### Issue 2 (ML Reproducibility Challenge 2020)

1. Replication in ML Reproducibility Challenge 2020 (Python) | 10.5281/zenodo.4835602 | PDF | Code | Review | BibTeX

VERMA, R., WAGEMANS, J.J.O., DAHAL, P., AND ELFRINK, A. 2021. [Re] Explaining Groups of Points in Low-Dimensional Representations. *ReScience C* 7, 2, #24.

2. Replication in ML Reproducibility Challenge 2020 (Python) | 10.5281/zenodo.4833219 | PDF | Code | Data | Review | BibTeX ALBANIS, G., ZIOULIS, N., CHATZITOFIS, A., DIMOU, A., ZARPALAS, D., AND DARAS, P. 2021. [Re] On end-toend 6DoF object pose estimation and robustness to object scale. *ReScience C* 7, 2, #2.

3. Replication in ML Reproducibility Challenge 2020 (python) | 10.5281/zenodo.4833389 | PDF | Code | Review | BibTeX ARVIND, M. AND MAMA, M. 2021. [Re] Neural Networks Fail to Learn Periodic Functions and How to Fix It. ReScience C 7, 2, #3.

#### RESCIENCE C

### IMPACT OF CHECKLISTS AND CHALLENGES

- Increase in the amount of code released during submission
- Increased interaction with authors and practitioners after paper publication through OpenReview





- Config management
- Logging
- Experimental Management
- Versioning
- Data management
- Data analysis
- Reporting
- Dependency Management
- Open Source Release
- Effective Communication
- Test and Release

Link to our previous blog post: <u>https://bit.ly/3LoSuKC</u>

- Config management
- Logging
- Experimental Management
- Versioning
- Data management
- Data analysis
- Reporting
- Dependency Management
- Open Source Release
- Effective Communication
- Test and Release



#### Hydra: https://hydra.cc

•••

Or even plain YAML / JSON files work! general: batch\_size: 128 data\_name: fashionmnist description: This is a sample config device: cuda epochs: 20 .ogbook: logger\_file\_path: log.jsonl log\_interval: 100 project\_name: fancy\_project nodel: class order: 0,1,2,3,4,5,6,7,8,9 loss\_policy: recon\_bce # ce, recon\_ce, recon\_mse, bce, recon\_bce max class: 10 reset optim: False eps: 1.0e-08 learning\_rate: 0.001 name: Adam scheduler gamma: 0.999 scheduler\_patience: 10 scheduler\_type: plateau weight decay: 0.0 sample mode: max vae\_hidden\_dim: 50 z dim: 5 resnet: in\_channels: 1

- Experimental Config management
- Logging •
- **Experimental Management** .
- Versioning •
- Data management ٠
- Data analysis
- Reporting ٠
- **Dependency Management** •
- **Open Source Release** •
- Effective Communication •
- Test and Release •



#### **Tensorboard**



STEP

Runs

#### Q Filter tags (regular expressions supported) Show data download links Ignore outliers in chart scaling epoch\_accuracy fooltip sorting method: default enoch accurac RELATIVE WALL 0 = 🖸 epoch\_loss 70100225-192554 (main 20190225-183554/validatio epoch\_loss 20190225-183652/data :: = ::

#### Weights & Biases

- Config management
- Logging
- Experimental Management
- Versioning
- Data management
- Data analysis
- Reporting
- Dependency Management
- Open Source Release
- Effective Communication
- Test and Release



Packaging format for

reproducible runs

on any platform

Record and query experiments: code,

data, config, results

General format for

sending models to

diverse deploy tools

- Config management
- Logging
- Experimental Management
- Versioning
- Data management
- Data analysis
- Reporting
- Dependency Management
- Open Source Release
- Effective Communication
- Test and Release

Add DeeBERT (entropy-based ea	arly exiting for *BERT) (#5477)	
🤔 Ji-Xin committed 10 days ago 🗸		
* Add deebert code		
* Add readme of deebert		
* Add test for deebert		
Update test for Deebert		
* Update DeeBert (README, clas	s names, function refactoring); remove	requirements.txt
* Format update		
* Update test		
* Update readme and model init	methods	
joeddav committed 10 days ago		
* add first draft ppl guide		
* upload imgs		
<pre>* expand on strides</pre>		
* ref typo		
* rm superfluous past var		Accelerating repro
* add tokenization disclaimer		
readme for benchmark (#5363)		1255
patrickvonplaten committed 10 d	lays ago 🗸	
mbart.prepare_translation_batch	h: pass through kwargs (#5581)	
Sshleifer committed 10 days ago	~	
Add mbart-large-cc25, support Sshleifer committed 10 days ago	translation finetuning (#5129)	
Create xlm-roberta-large-finetu	ned-conll03-german-README.md	8



lab

computational research.

- Config management
- Logging
- Experimental Management
- Versioning
- Data management
- Data analysis
- Reporting
- Dependency Management
- Open Source Release
- Effective Communication
- Test and Release





#### DVC, <u>https://dvc.org/</u>

#### **Datasheets for Datasets**

TIMNIT GEBRU, Google JAMIE MORGENSTERN, Georgia Institute of Technology BRIANA VECCHIONE, Cornell University JENNIFER WORTMAN VAUGHAN, Microsoft Research HANNA WALLACH, Microsoft Research HAL DAUMÉ III, Microsoft Research; University of Maryland KATE CRAWFORD, Microsoft Research; AI Now Institute

- Config management
- Logging
- Experimental Management
- Versioning
- Data management
- Data analysis
- Reporting
- Dependency Management
- Open Source Release
- Effective Communication
- Test and Release

Relevant works:

https://github.com/EleutherAl/Im-evaluation-h arness

	2' master + ParlAI / projects / controllable_dialogue / Analysis_n_Graphs.ipynb	Go to file	,	
	😮 stephenroller Release remaining Controllable Dialogue Code (#1734) 📖 🗸 Latest commit fb4b54d on Jun	2, 2019 🔇	History	,
	At 1 contributor			
	2.39 MB	Download	<b>₽</b> Ů	
	Evaluation Analysis (Public Release) Author: Stephen Roller <u>roller@tb.com</u> . Please direct questions to the ParAI Github issues ( <u>https://gthub.com/facebookresearch/ParAI/issues</u> ) This notebook expects to be iaunched from inside your ParAI installation (typically -/ParAI) You will need to pip install a bunch of things, including pyro and pandas. General preparation			
	<pre>In [33]: # bunch of imports and settings import os</pre>			
Jupyter	<pre>Specificity Control Level (WD) In [14]: def plot_resp.wd(matric, figgea, abslim, xaxis='Response-relatedness Control Level (WD)', use_title=False): plot_by_buckt(     modeltype_subset(altered, ["responsive"]),     modeling_subset(altered, ["respon</pre>			
	36- 14- 22-			



- Config management
- Logging
- Experimental Management
- Versioning
- Data management
- Data analysis
- Reporting
- Dependency Management
- Open Source Release
- Effective Communication
- Test and Release

#### The Machine Learning Reproducibility Checklist (v2.0, Apr.7 2020)

For all models and algorithms presented, check if you include:

- □ A clear description of the mathematical setting, algorithm, and/or model.
- A clear explanation of any assumptions.
- An analysis of the complexity (time, space, sample size) of any algorithm.

#### For any theoretical claim, check if you include:

- A clear statement of the claim.
- A complete proof of the claim.

For all datasets used, check if you include:

- □ The relevant statistics, such as number of examples.
- The details of train / validation / test splits.
- An explanation of any data that were excluded, and all pre-processing step.
- A link to a downloadable version of the dataset or simulation environment
- For new data collected, a complete description of the data collection process, such as instructions to annotators and methods for quality control.

For all shared code related to this work, check if you include:

- Specification of dependencies.
- Training code.
- Evaluation code.
- (Pre-)trained model(s).
- README file includes table of results accompanied by precise command to run to produce

#### **Model Cards for Model Reporting**

Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, Timnit Gebru {mmitchellai,simonewu,andrewzaldivar,parkerbarnes,lucyvasserman,benhutch,espitzer,tgebru}@google.com deborah.raji@mail.utoronto.ca
- Config management
- Logging
- Experimental Management
- Versioning
- Data management
- Data analysis
- Reporting
- Dependency Management
- Open Source Release
- Effective Communication
- Test and Release



- Config management
- Logging
- Experimental Management
- Versioning
- Data management
- Data analysis
- Reporting
- Dependency Management
- Open Source Release
- Effective Communication
- Test and Release



#### Language Models are Few-Shot Learners

28 May 2020 • Tom 8. Brown • Benjamin Mane • Neck Pyder • Malanis Subbiab • Jand Kaplan • Prafula Dharknal • Arvin Meekakantan • Prance Shyam • Ginish Sastry • Amanda Askell • Sandhini Agarwal • Arkil Hubert-Voss • Gentzben Kouger • Tom Hengham • Revon Child • Adhya Ramath • Daniel M. Zingler • Jaffrey Nu • Gamea Winter • Articother Henses • Mark Chen • Eric Sigler • Matexa: Univer • Sont Gray • Benjimin Charol • Jaka • Maya Ramath • Daniel M. Zingler • Jaffrey Nu • Gamea Winter • Andre Matexa • Leic Sigler • Matexa: Univer • Sont Gray • Benjimin Charol • Jaka • Bahar • Maya Ramath • Charolin • Jaka • Radinat • Ng Sastware • Danie Andrei

Recent work has demonstrated substantial gains on many NLP tasks and benchmarks by pre-training on a large corpus of text followed by fine-tuning on a specific task. While typically task-agnostic in architecture, this method still requires task-specific fine-tuning datasets of thousands or tens of thousands of examples... (read more)



Code		🖉 Edit	Tasks	🖪 Edit
O openai/gpt-3	★ 5,107		COMMON SENSE REASONING	
O sw-yx/gpt3-list	★ 95		COREFERENCE RESOLUTION	
() Tacebookresearch/anii	★ 83		DOMAIN ADAPTATION	
			FEW-SHOT LEARNING	

- Config management
- Logging
- Experimental Management
- Versioning
- Data management
- Data analysis
- Reporting
- Dependency Management
- Open Source Release
- Effective Communication
- Test and Release



NeurIPS 2019 repositories with 0 ticks had a median of 1.5 GitHub stars. In contrast, repositories with 5 ticks had a median of 196.5 GitHub stars. Only 9% of repositories had 5 ticks, and most repositories (70%) had 3 ticks or less.

- Config management
- Logging
- Experimental Management
- Versioning
- Data management
- Data analysis
- Reporting
- Dependency Management
- Open Source Release
- Effective Communication
- Test and Release





Google Colab





**Google** Cloud Platform



CO CML by iterative.ai

# QUESTIONS?

113

# MOTIVATION

How can we mitigate the challenges of bigger, more complex models without reducing the benefits?

# In this tutorial, we focus on the challenge of <u>ensuring</u> <u>research results are reproducible</u>

# TUTORIAL OVERVIEW

- I. Introduction to Reproducibility
- 2. Reproducibility in IR
- 3. Mechanisms for Reproducibility
- 4. Reproducibility as a Teaching Tool

# IV – REPRODUCIBILITY AS A TEACHING TOOL

Maurits Bleeker, Ana Lucic

# OVERVIEW

- 1. Teaching through reproducibility
- 2. Examples of AI courses utilizing reproducibility as a teaching tool
  - a. Reproducedpapers.org (TU Delft)
  - b. FACT-AI (University of Amsterdam)
- 3. Guidelines for a successful reproducibility course
- 4. Lessons learned

# IV.a – TEACHING THROUGH REPRODUCIBILITY

# EXAMPLES FROM OTHER ACADEMIC FIELDS

- Learning Networking by Reproducing Research Results (Yan et al. 2017)
  - Stanford CS course on reproducing work on networking systems
- Bringing Replication Into Classroom: Benefits For Education, Science, and Society (Ribotta, Blandine, et al 2022)
  - "For more than a decade, research in psychology has been struggling to replicate many well-known and highly cited studies"
- How to Use Replication Assignments for Teaching Integrity in Empirical Archaeology (Marwick, Ben, et al. 2020)
  - "Here we argue for replications as a core type of class assignment in archaeology courses"

# MOTIVATION

Valuable experience for students:

- Practice implementing and extending existing research
- Recognize the importance (and difficulty) of reproducibility
- Helps students to develop critical thinking skills
  - This also helps with writing research papers
- Can be added to their portfolio, e.g., personal website, blog post, CV
  - Allows students to participate in the community

Contribute to existing research:

- New insights can direct future research
- Results can be published, e.g., in the *ReScience journal*

# IV.b - REPRODUCEPAPERS.ORG

# REPRODUCEDPAPERS.ORG

"Is an open online repository for teaching and structuring machine learning reproducibility"

- Primary motivation: there exist several venues for reproducibility but there is a 'high barrier' to entry or a focus on 'short-term' (alternate years, etc)
- Propose: a low barrier, long term venue focused on reproducibility
- Reproduction aligns with several teaching goals:
  - Reading and critiquing literature
  - Implementing, executing and extending code
  - Comparing, analyzing and presenting results in a clear and concise manner

# ONLINE REPOSITORY

Search for papers and reproductions	Reproductions Papers Help About <b>TU</b> Delft Sig
Reproductions	Submit Reproduction
Reproduced New data New algorithm variant Reproduction of "SwinIR: Image Restorated by Frans de Boer, Jonathan Borg, Adarsh Denga, Ha We explain the technical details of the SwinIR paper in our own word algorithm. Furthermore we explore modifying the architecture used energy. Detail	tion Using Swin Transformer" aoran Xia Is, providing ample detail to understand the authors' contribution and in the paper to allow it to run using reduced resources and thus use less
Replicated Reproduction of "Deep Learning with Diff by Deep Learning CS4240 Group66: Hengkai Zhan Benefits of machine learning techniques based on neuron networks sensitive information should be retained. Differential privacy is thus "Deep More	<b>ferential Privacy"</b> Ig, Dong Shen, Yuxin Cheng are widely appreciated. While these methods require a large amount of data, developed. This blog aims to present and describe our efforts to reproduce

# **ONLINE REPOSITORY**

- Focus of the project: partial results, minor tweaks, etc.
- Well suited for use in teaching
- Badges (self-labeled):
  - **Replicated:** A full implementation from scratch without using any pre-existing code
  - **Reproduced:** Existing code was evaluated
  - **Hyperparams check:** New evaluation of hyperparameter sensitivity
  - New data: Evaluating new datasets to obtain similar results
  - **New algorithm variant:** Evaluating a different variant
  - **New code variant:** Rewrote/ported existing code to be more efficient /readable
  - Ablation study: Additional ablation studies

# COURSE DETAILS

- Part of MSc CS Deep Learning course, TU Delft
- Teaching team selects papers with two criteria:
  - Data availability
  - Computational demands
- Projects:
  - Teams should indicate which result to reproduce
  - Groups of 2-4 students, 8 week course
  - $\frac{1}{3}$  of the course time spent on reproduction
- Deliverables:
  - Blog about the repository (private/public)
  - PDF report

24 unique papers, 57 paper reproductions



Yildiz et al. 2021. ReproducedPapers.org: Openly teaching and structuring machine learning reproducibility. International Workshop on Reproducible Research in Pattern Recognition.



Yildiz et al. 2021. ReproducedPapers.org: Openly teaching and structuring machine learning reproducibility. International Workshop on Reproducible Research in Pattern Recognition.



Yildiz et al. 2021. ReproducedPapers.org: Openly teaching and structuring machine learning reproducibility. International Workshop on Reproducible Research in Pattern Recognition.

# CONCLUSION

- Reproduction projects align closely with general course learning goals, and were received positively by most students
- These projects improve perceived value of reproductions, with an added incentive of publishing their work and adding to their portfolios
- "We finally call on the community to add their reproductions to the website ReproducedPapers.org"
- "May the next generation of machine learners be reproducers"

## IV.c – FAIRNESS, ACCOUNTABILITY, CONFIDENTIALITY, AND TRANSPARENCY IN AI COURSE

University of Amsterdam

# COURSE MOTIVATION

- In 2019, we designed a new course on Fairness, Accountability, Confidentiality, and Transparency in AI (FACT-AI) at the University of Amsterdam (UvA)
  - Based on the requests of our students in the MSc AI: an increase in interest in ethical issues in AI
- The course aims to make students aware of two types of responsibility:
  - Towards society in terms of potential implications of their research
    - Similar to the NeurIPS Paper Checklist: discuss any potential negative societal impacts of your work
  - Towards the research community in terms of producing reproducible research



- LO #I: Understanding FACT topics
- LO #2: Understanding algorithmic harm
- LO #3: Familiarity with FACT methods
- LO #4: Reproducing FACT solutions

Learning Objective #1: Understanding FACT topics

• Students can explain the major notions of fairness, accountability, confidentiality, and transparency that have been proposed in the literature, along with their strengths and weaknesses

## Learning Mechanism:

• General lecture(s) per topic

**Learning Objective #2:** Understanding algorithmic harm

• Students can explain, motivate, and distinguish the main types of algorithmic harm, both in general and in terms of concrete examples where AI is being applied

## Learning Mechanism:

- General lectures and guest lectures, where students can ask questions and are encouraged to participate in discussions
- This LO can be used for any AI course

Learning Objective #3: Familiarity with FACT methods

• Students are familiar with recent peer-reviewed algorithmic approaches in the FACT-AI literature

## Learning Mechanism:

• Paper discussion sessions where students discuss a seminal FACT-AI paper in a small and interactive group, after reading the paper in advance

# PAPER DISCUSSION

- Outline of how to dissect a paper ahead of time
  - Examples help!
- For the students, the goal of the paper discussion sessions is to:
  - Learn about prominent methods in the field
  - Reading a technical paper
  - Think critically about the claims made in the papers
  - Understanding a paper's strength and weaknesses
- All these (reading) skills are necessary for a good reproducibility study
  - If students can't understand the paper, how will they reimplement the algorithm?

# PAPER DISCUSSION

- Students first read a seminal paper on their own trying to answer the following questions:
  - What are the main claims of the paper?
  - What are the research questions?
  - Does the experimental setup make sense, given the research questions?
  - What are the answers to the research questions? Are these supported by experimental evidence?
- Participate in small discussion sessions (ideally in person) with their peers to discuss their answers
  - Groups of 4 to 5 students

# PAPER DISCUSSION

An instructor goes over the same paper, giving an overview of the papers' strengths and weaknesses

- In our case, each session was presented by a different instructor
- This to show:
  - There is no single way of examining a research paper
  - Different researchers will bring different perspectives to their assessment of papers
- We chose papers for their discussion sessions based on their impact on the FACT-AI field

Learning Objective #4: Reproducing FACT solutions

• Students can assess the degree to which recent algorithmic solutions are effective, especially with respect to the claims made in the original papers, while understanding their limitations and shortcomings

## Learning Mechanism:

Group project where students work in groups to reproduce FACT-AI papers from top AI conferences

- The group project is based on **reproducing existing algorithms** from top AI conferences and is the focal point of the course
- In our course, we focused on FACT-AI algorithms
- However, the setup for the course is not specific to FACT-AI and can be tailored to any topic
  - e.g., **IR**, computer vision, information retrieval, general ML, etc.

- The group project is based on **reproducing existing algorithms** from top AI conferences and is the focal point of the course
- Students work in groups to reimplement existing algorithms from papers in top Al conferences (e.g., NeurIPS, ICML, ICLR, AAAI, etc).
- Students write up the results and submit reports
  - We encouraged them to submit their reports to the ML Reproducibility Challenge
- In our course, we focused on FACT-AI algorithms. However, the setup for the course is not specific to FACT-AI and can be tailored to any topic
  - e.g., , computer vision, **information retrieva**l, general ML, etc.

Benefits of participating in the ML Reproducibility Challenge:

- Motivates and incentivizes students
- Reports accepted by the ML Reproducibility Challenge are accepted for publication in the *ReScience* journal
- Exposes students to the paper submission cycle

Participating in the ML Reproducibility Challenge gives the students the opportunity to experience the whole research pipeline:

- 1. Reading a technical paper to understand its strength and weaknesses
- 2. Implementing (and perhaps also extending) the algorithms in the paper
- 3. Writing up the findings
- 4. Submitting to a venue with a deadline
- 5. Obtaining feedback from reviewers
- 6. Writing a rebuttal
- 7. Receiving the official acceptance/rejection notification
## COURSES PARTICIPATE IN RC2021 FALL EDITION

#### Courses Participated in RC2021 Fall Edition

- DD2412 Deep Learning, Advanced. KTH (Royal Institute of Technology), Stockholm, Sweden
- CISC 867 Deep Learning, Queen's University, Ontario, Canada
- Special Topics in CSE: Advanced ML, Indian Institute of Technology, Gandhinagar, India
- FACT: Fairness, Accountability, Confidentiality and Transparency in AI, University of Amsterdam,

#### Netherlands

- CSCI 662 -- Advanced Natural Language Processing, University of Southern California, USA
- Intelligent Systems and Interfaces, Indian Institute of Technology, Guwahati, India
- Intelligent Information Processing Topics, Tsinghua University, China
- Machine learning for data science 2, University of Ljubljana, Slovenia
- EECS 598-005: Randomized Numerical Linear Algebra in Machine Learning, University of Michigan, USA
- SYDE 671 Advanced Image Processing, University of Waterloo, Canada
- BLG561E Deep Learning, Istanbul Technical University, Turkey
- CS 433 Machine Learning, EPFL, Switzerland

# **RESULTS OF THE ML REPRODUCIBILITY CHALLENGE**

- See <a href="https://openreview.net/group?id=ML\_Reproducibility\_Challenge">https://openreview.net/group?id=ML\_Reproducibility\_Challenge</a>
- ML Reproducibility Challenge 2021
  - ± 40% of the accepted papers were from the UvA FACT-AI course
- ML Reproducibility Challenge 2022
  - ± 50% of the accepted papers were from the UvA FACT-AI course
  - Best paper award
  - 2 outstanding papers (out of 4)

#### FEEDBACK

First year MSc AI students

"I appreciate the critical view I have developed on papers as a result of this course. Normally I would easily accept the content of a paper, but I will be more critical from now on, as many papers are not reproducible."

"I really appreciated that this was the first course where students are judging state-of-the-art AI models. In other words, students were able to experience the scientific workfield of AI."

#### FEEDBACK

First year MSc AI students

"Replicating another study, seeing how (poorly) other research is performed was really eye-opening."

"I think it's really good that we get some practical insights into reproducing results from other papers, not all papers are as good as they seem to be."

# QUESTIONS?

## IV.d – GUIDELINES FOR A SUCCESSFUL REPRODUCIBILITY COURSE

#### GUIDELINES FOR A SUCCESSFUL REPRODUCIBILITY COURSE

- INCLUDE A REPRODUCIBILITY LECTURE
- PAPER REQUIREMENTS
- GRADING
- TEACHING ASSISTANTS
- TIMING OF THE COURSE
- DURATION OF THE COURSE
- ADVANTAGES OF PARTICIPATING IN THE ML REPRODUCIBILITY CHALLENGE

# INCLUDE A REPRODUCIBILITY LECTURE

Motivate reproducibility with a general lecture

- Position this lecture (ideally) at the beginning of the course
- Highlight papers examining reproducibility/replicability failures
  - For examples in IR, see Part 2 of the tutorial
  - Include consequences of failure to reproduce (Part 2)
- Clearly outline scope of the project(s) and potential impact

# PAPER REQUIREMENTS

- Choose 10-15 papers from the ML Reproducibility Challenge OpenReview portal that are suitable for your course
- Before the course starts, let the TAs check whether the selected papers are feasible for reproducibility study
  - Hire a team of experienced, graduate-level TAs
- Ideally assign each TA no more than 3-4 papers

# PAPER REQUIREMENTS

- Select papers that are computationally feasible to reproduce
  - In our case, we were able to provide one GPU per team
  - Depends the available resources of the course and faculty
- At least one dataset should be publicly available and of a reasonable size
  - If the dataset is too big, it is an option to reproduce the work in a 'low-resource' data setting
- Select papers that are relevant to the topics covered in the course
- Emphasize the technical perspective of the sub-field
- It should be reasonable to reimplement the paper within the allotted time

# GRADING

- Grading group projects on different papers in a fair manner is challenging
- Try to make the grading criteria as explicit as possible in order to make it clear for the students what is expected
- Organize a grade calibration session with the TAs after grading to align on expectations
- If participating in the ML Reproducibility Challenge, grade reports independently of the reviews

		Grade	<= 5 (fail)	6 (sufficient)	7 (satisfactory)	8 (good)	9 (very good)	10 (excellent)
	Project (40%)							
		Project Design	Unsystematic and/or no validated use of research and design methodologies. Insufficient explanation. How are the results tested and/or verified?	Adequate use of research and design methodologies. Limited explanation.	Adequate use of research and design methodologies. Explained and justified.	Use of the right research and design methodologies. Well-explained and well justified.	Profound and critical use of research and design methodologies. Very clear and validated design.	Excellent demonstration of research and design methodologies.
		Positioning of project	Project not positioned w.r.t. new literature, the FACT-field and reproducibility papers.	Project is somewhat positioned.	Project is sufficiently positioned in literature.	Project is correctly positioned in literature.	Project is well positioned within literature.	Project is integrated within literature, even from different fields/sources.
		Creativity	The project does not make an original contribution. E.g. the picked paper is just said to be reproducible or not without any extra insights.	Project does not really make any original contribution. The results are reproducible, with limited effort or not reproducible with limited insights (why is this not working?).	Project team had at least one original contribution to reproduce the work and/or go beyond the original results of the paper.	Project team came up with several original ideas to reproduce the paper and/or go beyond the original results, design options and/or concepts not initiated or thought of by the supervisor.	Project team came up with many original ideas, design options and/or concepts to reproduce the work and/or go beyond the originial results. Not initiated or thought of by the supervisor.	Project team surprised us all with some brilliant new ideas, design options and/or concepts, both in breadth and depth.

Code base (20%)								
	Technical quality	Insufficient	Sufficient	Satisfactory	Good	Very Good	Excellent	
	Reproducability of your results by the TA's.	Not reproducible. The project results should be reproducible by the TAs	N/A	With some effort the results are reproducible by the TAs.	N/A	Without any effort the results are reproducible by the TAs	N/A	

Paper (30%)										
		Content	Report shows no coherence of content. For example: What questions are you asking? What experiments do you run to answer them? What conclusions can you draw from these experiments?	Report shows sufficient coherence of content.	Report fulfils all requirements in terms of content.	Good report in terms of content.	Very good report in terms of content.	Excellent report in terms of content.		
		Form	Structure needs considerable improvement. General presentation of the content (text and figures) not very effective.	Structure needs some improvement. General presentation of the content (text and figures) is sufficient.	Structure is acceptable. General presentation of the content (text and figures) is satisfactory.	Clear structure. Good presentation of the content (text and figures).	Well-structured document. General presentation of the content (text and figures) is effective.	Very well-structured document. General presentation of the content (text and figures) is very effective.		
		Quality of writing	Poorly expressed. Document contains serious spelling and grammatical errors.	Reasonably expressed argumentation. Document contains some spelling and grammatical errors.	Sufficiently expressed argumentation. The document contains little spelling and grammatical errors.	Expressed and formulated well. Document has a nice flow. Document contains only minor spelling and grammatical errors.	Expressed and formulated very well. Document has a smooth flow with sufficient transitions. Document is without any spelling and grammatical errors.	Excellent expressed and formulated report. Document has a smooth flow with effective transitions. Spelling and grammatically error free.		

Presentation (10%)							
	Content	Presentation lacks detail and does not support conclusions. Irrelevant information presented.	Presentation lacks detail, and is just enough to support conclusions.	Presentation has sufficient detail to support conclusions.	Presentation has a good level of detail to support conclusions.	Presentation has the right level of detail to support the conclusions and to understand the recommendations.	Presentation has the perfect level of detail to support the conclusions and to understand the recommendations.
	Form	Presentation is unstructured and not well organized. No (proper) use of visual aids.	Logical structure of presentation is poor. Improvements to the structure should be made. Use of visual aids can be improved.	Logical structure of presentation is reasonable but needs some improvement. Sufficient use of visual aids.	Presentation has good logical structure, the essentials are separated from the ancillary. Good use of visual aids.	Presentation has very good logical structure, the essentials are clearly separated from the ancillary. Good use of visual aids.	Presentation has excellent logical structure, the essentials are very well separated from the ancillary. Perfect use of visual aids.
	Performance	Poorly expressed and formulated. Unclearly presented. Audience was ineffectively addressed.	Expression and formulation can be improved. Not always clearly presented.	Expressed and formulated adequately. Most of the time clearly presented. Audience was sufficiently addressed.	Well expressed and formulated. Clearly presented. Audience was well addressed.	Very well expressed, formulated and clearly presented.	Expressed, formulated and presented with great style, clarity and effectiveness. Audience was very well addressed and engaged.

# **TEACHING ASSISTANTS**

- Have the TAs read the papers before the course starts to ensure they have a sufficient, in-depth understanding of their papers
  - Assign papers to TAs based on their interests
- To ease the load for the TAs, have several groups working on the same paper
- Ensure students have regular contact with their TA so no group gets stuck in the process
- Ask students halfway through the course to submit a draft report to their TAs in order to get feedback
  - We found this significantly increased the quality of the final reports

# TIMING OF THE COURSE

#### Students need to have very strong programming skills

Table 1: The first year of the MSc AI program at the University of Amsterdam.						
Course	Sem. 1	Sem. 2	EC			
Computer Vision 1			6			
Machine Learning 1			6			
Natural Language Processing 1			6			
Deep Learning 1			6			
Fairness, Accountability, Confidentiality			6			
and Transparency in AI						
Information Retrieval 1			6			
Knowledge Representation and Reasoning			6			
Elective 1			6			
Elective 2			6			
Elective 3			6			

# DURATION OF THE COURSE

- We strongly recommend to ensure that the students to have enough time to work on the project
- For our course, the students are working one month full-time on the project
  - We found this to be a beneficial setup since students didn't have to worry about any other courses during this time
- If it's not possible to work on the project full-time, then potentially adapt the weight of the course:
  - If students typically have 5 courses in one semester, consider making the reproducibility course worth 2 courses

# ADVANTAGES OF PARTICIPATING IN THE ML REPRODUCIBILITY CHALLENGE

- Prioritize the ML Reproducibility Challenge by tying the reproducibility report directly to the grading
  - Students are graded on the same report that they submitted to the challenge therefore, participating is not an extra task
- Submitting to the challenge gives the students the opportunity to experience the whole research pipeline:
  - Submitting to a venue with a strict deadline
  - Obtaining feedback
  - Writing a rebuttal
  - Receiving the official notification

## IV.e – LESSONS LEARNED

# SUMMARY OF THE LESSONS LEARNED

In our experiences, we found that the following were important components of a successful course:

- Including extension as part of reproducibility
- Having excellent teaching assistants
- Having students participate in the ML community
- Encouraging communication with the original authors

# INCLUDING EXTENSIONS AS PART OF REPRODUCIBILITY

- We argue that the finding *"the original work is (not) reproducible"* is not insightful
- Require students to extend the paper if the source-code is already available
- Either extend the work to:
  - New domains, datasets or a low-resource regime (i.e., less data/compute)
  - New hyper-parameter settings or method different assumptions
  - Different model architecture
- Or explain why the work is not reproducible

# INCLUDING EXTENSIONS AS PART OF REPRODUCIBILITY

There are two scenarios possible for the project:

- There already exists an open-source implementation of the selected paper. Students are allowed to use this:
  - The results the students obtain are different as described in the paper
  - The results are reproducible, meaning this method can now be used for further research
- There is no open-source implementation available, meaning the students need to reimplement everything themselves
  - Take this into account when grading

# HAVING EXCELLENT TEACHING ASSISTANTS

- It is extremely important for the TAs to have **excellent programming experience** since this is the main aspect students need help with
- Have students meet with the TAs at least twice a week
- We had both second year MSc students and PhD students
  - PhD students are prefered, if possible
- Have the TAs help students with writing the rebuttal, since this is a new experience for them

# HAVING EXCELLENT TEACHING ASSISTANTS

Since this is probably the first time the students are submitting a research paper, try to prevent the following common mistakes:

- Submitting single blind
- Referring to the course project in the introduction
- Motivation: "We had to do this for a course project"
- Submitting a non-anonymized code-base

#### HAVING STUDENTS PARTICIPATE IN THE ML COMMUNITY

- It is a motivating factor for students to create concrete output that is beneficial to the broader ML research community
- FACT-AI course 2019--2020
  - Creating a public repository with the best algorithm implementations
- FACT-AI course 2020--2021 and 2021--2022:
  - Participating in the ML Reproducibility Challenge

# ENCOURAGING COMMUNICATION WITH THE ORIGINAL AUTHORS

- We strongly encourage students to contact the original authors
- It is beneficial for students to interact with scientists in the field
- It improves the papers' credibility, readability, and reproducibility
- Give the students some instructions how to do this:
  - Be aware that the authors are busy
  - Prevent that multiple teams are emailing at the same time
    - Have the TAs coordinate this

# SUMMARY OF THE LESSONS LEARNED

In our experiences, we found that the following were important components of a successful course:

- Including extension as part of reproducibility
- Having excellent teaching assistants
- Having students participate in the ML community
- Encouraging communication with the original authors

# QUESTIONS?

# V. – CONCLUSION

# CONCLUSION

- We have shown two successful examples of graduate-level AI courses that focus on reproducibility with their course project
- We provided guidelines to successfully run a reproducibility project for any graduate-level AI course
- Implementing a course centred on a reproducibility project is fairly straightforward for the instructor and has many benefits for students
  - The course naturally "refreshes" itself every year when a new batch of papers is chosen

# MOTIVATION

How can we mitigate the challenges of bigger, more complex models without reducing the benefits?

# In this tutorial, we focus on the challenge of <u>ensuring</u> <u>research results are reproducible</u>

# SUMMARY OF TUTORIAL

In this tutorial, we've aimed to address the issue of **ensuring research results are reproducible** 

- Part 1:We gave an introduction to reproducibility and presented some examples of (ir)reproducible results, both from within CS and from other disciplines
- Part 2:We went over reproducibility aspects in IR as well as some examples of reproducibility failures and ongoing efforts to help improve reproducibility in IR
- Part 3:We investigated existing mechanisms for reproducibility in ML/IR such as Papers with Code and the ML Reproducibility Challenge
- Part 4:We discuss how to teach reproducibility to the next generation of AI researchers

# BEST PRACTICES TO KEEP IN MIND

- 1. **Report** as much as much information as you can
  - Different types of papers have different requirements when creating a new dataset, consider the annotators! When running experiments, do a hyperparameter search!
- 2. Share dependency config files
- 3. Release code
  - If an experiment didn't work or provides evidence that doesn't support your main hypothesis (e.g., that your model is better than previous models), you should still report it!
- 4. **Run** multiple experiments (with different random seeds, or different data orders, etc.) and report error bars.
- 5. **Record** your carbon emissions
  - You can use tools like <u>CodeCarbon</u> or the <u>ML CO2 Calculator</u>
- 6. **Fill out** reproducibility checklists correctly, try to do any items that are appropriate (though we recognize the checklists aren't perfect)