

MAKING YOUR RESEARCH REPRODUCIBLE Odd Erik Gundersen, dr. philos.

Chief AI Officer, TrønderEnergi AS Adjunct Associate Professor, NTNU



Norwegian Open Al Lab





AlphaGo



"Impressive results. No code. No model."







"I think you should be more explicit here in step two."



Reproducing AlphaZero with ELF: What we learned

Yuandong Tian Facebook AI Research



Yuandong Tian



Jerry Ma*







Qucheng Gong* Shubho Sengupta* Zhuoyuan Chen James Pinkerton





Larry Zitnick







Reproducing AlphaZero with Elf

- Hard to reproduce \bullet
 - **Details are missing in the paper** _

 - Sophisticated (distributed) systems.
- Lack of ablation analysis \bullet
 - What factor is critical for the performance? _

 - How the ladder issue is solved?
- Lots of mysteries \bullet
 - Is the proposed algorithm really universal? _____
 - Is the bot almighty?
 - Is there any weakness in the trained bot?

ELF OpenGo: An Analysis and Open Reimplementation of AlphaZero, Tian et al, ICML 2019.

Huge computational cost (15.5 years to generate 4.9M selfplays with 1 GPU)

Is the algorithm robust to random initialization and changes of hyper parameters?



52% Yes, a significant crisis

IS THERE A REPRODUCIBILITY CRISIS?

A Nature survey lifts the lid on how researchers view the 'crisis' rocking science and what they think will help.

1,576 **RESEARCHERS SURVEYED**

7% Don't know

3% No, there is no crisis

BY MONYA BAKER

38% Yes, a slight crisis

(M. Baker, Nature, 2016)

(Gundersen , 2020)



HAVE YOU FAILED TO REPRODUCE AN EXPERIMENT? Most scientists have experienced failure to reproduce results. My own Someone else's ÷. . . Chemistry 44 1. 50 Biology Physics and engineering 0 Medicine ÷. 2.1 Earth and environment Other 20 40 60 80 0



Computer Science





WHAT FACTORS CONTRIBUTE TO **IRREPRODUCIBLE RESEARCH?**

Many top-rated factors relate to intense competition and time pressure.

Always/often contribute
Sometimes contribute

Selective report

Pressure to pul

Low statistical power or poor ana

Not replicated enough in origina

Insufficient oversight/mento

Methods, code unavail

Poor experimental de

Raw data not available from origina

Insufficient peer re

rting						
_						
blish						
-						
alysis						
-						
al lab						
oring						
-						
lable						
esign						
-						
al lab						
_						
raud						
_						
eview						
///0//						
()	20	40	60	80	100%





MISSING DATA

As research articles age, the odds of their raw data being extant drop dramatically.





ICLR 2018 Reproducibility Challenge

Before the challenge (n=98): "Is there a reproducibility crisis in ML?"





(J. Pineau, ICLR keynote, 2018)





REPRODUCIBILITY PART II





The Scientific Method - Process





The Scientific Method - Steps

- Observe the world and form beliefs about it 1.
- Explain causes and effects by forming a scientific theory 2.
- 3. Formulate a genuine test of the scientific theory as a hypothesis
- Design an experiment to test the hypothesis and document it in a research 4. protocol
- Implement the experiment so that it is ready to be conducted
- 5. 6. Conduct the experiment to produce results
- 7. Analyze the results to make an analysis
- 8. Interpret the findings
- Update beliefs according to the interpretation 9.
- Observe the world in a structured manner 10.





Types of Empirical Studies

Hypothesis generating



Hypothesis generating - identify and suggest possible hypotheses.

Hypothesis testing - test explicit and precise hypotheses

- \bullet measure variables.





Exploratory Yields casual hypotheses by collecting data and analyzing it in many ways.

Assessment Establish baselines and ranges as well as other behaviors of system or environment.

Observation Collect data in a way that does not directly interfere with how the data arise, establish an association. *Manipulation* Test hypotheses about causal influences of factors by manipulating them and and noting effects on





The Scientific Method in ML









Example of Experiment

Multi-column Deep Neural Networks for Image Classification

Dan Cireşan, Ueli Meier and Jürgen Schmidhuber IDSIA-USI-SUPSI Galleria 2, 6928 Manno-Lugano, Switzerland

{dan,ueli,juergen}@idsia.ch

Abstract

Traditional methods of computer vision and machine learning cannot match human performance on tasks such as the recognition of handwritten digits or traffic signs. Our biologically plausible, wide and deep artificial neural network architectures can. Small (often minimal) receptive fields of convolutional winner-take-all neurons yield large network depth, resulting in roughly as many sparsely connected neural layers as found in mammals between retina and visual cortex. Only winner neurons are trained. Several deep neural columns become experts on inputs preprocessed in different ways: their predictions are averaged. Graphics cards allow for fast training. On the very competitive MNIST handwriting benchmark, our method is the first to achieve near-human performance. On a traffic sign recognition benchmark it outperforms humans by a factor of two. We also improve the state-of-the-art on a plethora of common image classification benchmarks.

1. Introduction

Recent publications suggest that unsupervised pretraining of deep, hierarchical neural networks improves supervised pattern classification [2, 10]. Here we train such nets by simple online back-propagation, setting new, greatly improved records on MNIST [19], Latin letters [13], Chinese characters [22], traffic signs [33], NORB (jittered, cluttered) [20] and CIFAR10 [17] benchmarks.

We focus on deep convolutional neural networks (DNN), introduced by [11], improved by [19], refined and simplified by [1, 32, 7]. Lately, DNN proved their mettle on data sets ranging from handwritten digits (MNIST) [5, 7], handwritten characters [6] to 3D toys (NORB) and faces [34]. DNNs fully unfold their potential when they are wide (many maps per layer) and deep (many layers) [7]. But training them requires weeks, months, even years on CPUs. High data transfer latency prevents multi-threading and multi-CPU code from saving the situation. In recent years, however, fast parallel neural net code for graphics cards (GPUs)

has overcome this problem. Carefully designed GPU code for image classification can be up to two orders of magnitude faster than its CPU counterpart [35, 34]. Hence, to train huge DNN in hours or days, we implement them on GPU, building upon the work of [5, 7]. The training algorithm is fully online, i.e. weight updates occur after each error back-propagation step. We will show that properly trained wide and deep DNNs can outperform all previous methods, and demonstrate that unsupervised initialization/pretraining is not necessary (although we don't deny that it might help sometimes, especially for datasets with few samples per class). We also show how combining several DNN columns into a Multi-column DNN (MCDNN) further decreases the error rate by 30-40%.

2. Architecture

The initially random weights of the DNN are iteratively trained to minimize the classification error on a set of labeled training images; generalization performance is then tested on a separate set of test images. Our architecture does this by combining several techniques in a novel way:

(1) Unlike the small NN used in many applications, which were either shallow [32] or had few maps per layer (LeNet7, [20]), ours are deep and have hundreds of maps per layer, inspired by the Neocognitron [11], with many (6-10) layers of non-linear neurons stacked on top of each other, comparable to the number of layers found between retina and visual cortex of macaque monkeys [3].

(2) It was shown [14] that such multi-layered DNN are hard to train by standard gradient descent [36, 18, 28], the method of choice from a mathematical/algorithmic point of view. Today's computers, however, are fast enough for this, more than 60000 times faster than those of the early 90s¹. Carefully designed code for massively parallel graphics processing units (GPUs normally used for video games) allows for gaining an additional speedup factor of 50-100 over serial code for standard computers. Given enough labeled data, our networks do not need additional heuristics

Scientific theory: Deep neural networks are models of the brain, although simple ones, and as such intelligence could emerge from them.

Hypothesis: The performance of biological inspired deep convolutional neural networks is competitive with human performance on computer vision benchmark tasks.



^{11991 486}DX-33 MHz, 2011 i7-990X 3.46 GHz



Experiment: DNN for Image Classification







The Scientific Method - Steps







The Scientific Method in ML











Definition of Reproducibility

Reproducibility is the ability of independent investigators to draw the same conclusions from an experiment by following the documentation shared by the original investigators.







The Three Types of Reproducibility

Outcome reproducible The outcome of the reproducibility experiment is

the same as the outcome produced by the original experiment.

Analysis reproducible Outcome might differ, but same analysis and interpretation on different outcome leads to same conclusion.

Interpretation reproducible Neither the outcome nor the analysis need

to be the same if the interpretation leads to the same conclusion.





(Gundersen 2021)



The Three Types of Documentation

Description Description of the AI method implemented by the AI program, the

experiment being conducted and the analysis of the results as well as the

hardware and ancillary software used for conducting the experiment.

Code AI Program code, code for setup and configuration, code controling workflow,

code for analysis of results and visualization.

Data All data used for conducting the experiment. Are the samples used for

training, validation and test specified? What about the results?









The Scientific Method in ML







Result Experiment Outcome Analysis Interpretation

Degrees of Reproducibility

			1
	Outcome	Analysis	Interpretation
R1 Outcome	Same	Same	Same
R2 Analysis	Different	Same	Same
R3 Interpretation	Different	Different	Same
orld $\xrightarrow{1.}$ Beliefs $\xrightarrow{2.}$ Scientific theor	ry $\xrightarrow{3.}$ Hypothesis $\xrightarrow{4.}$ Rese	9. earch protocol $\xrightarrow{5.}$ Experiment	ment $\stackrel{6.}{\longrightarrow}$ Outcome $\stackrel{7.}{\longrightarrow}$ Analy
10. Explorat Assessm	tory studies Hypothesis - ent studies Prediction - Study plan - Analysis plan	Target system Experiment workflow Experiment setup Pre-processing code	Code for analysis Code for visualization
10).	Data Ancillary software Hardware	

undersen 2021)

Four Types of Reproducibility

	Text
Description	
Code	
Data	
Experiment	

Claerbout 1992, reproducibility -Buckheit 1995, reproducibility -Peng 2006, replication -Peng 2006, reproducibility -Drummond 2009, repoducibility -Drummond 2009, replicability -Schmidt 2009, direct replication -Schmidt 2009, conceptual replication -Miller 2010, repeatability -Miller 2010, reproducibility -Peng 2011, publication only -Peng 2011, full replication -Stodden 2011, replicability -Stodden 2011, reproducibility -JCGM 2012, repeatability -JCGM 2012, reproducibility -Crook 2013, internal replicability -Crook 2013, external replicatbility -Crook 2013, cross-replicability -Crook 2013, reproducibility -Gent 2014, recomputation -Nosek 2014, direct replication -Nosek 2014, conceptual replicaton -Goodman 2016, results reproducible -Goodman 2016, methods reproducible -ACM 2018, repeatibility -ACM 2018, replicability -ACM 2018, reproducibility -Gundersen 2018, experiment reproducible -Gundersen 2018, data reproducible -Gundersen 2018, method reproducible -NAS 2019, reproducibility -NAS 2019, replicability -

THE AAAI REPRODUCIBILITY CHECKLIST PART III

NTNU

General Reproducibility Guidelines for AI Research

Version: 1.3 June 25, 2020 Authors: Odd Erik Gundersen, Yolanda Gil, Mausam

For each experiment, check that the following is described:

- How the experimental design rigorously tests the claims.
- The evaluation metrics and the motivation for choosing these metrics.
- All (hyper-)parameters for each model/algorithm, number and range of values tried per parameter, and the criterion for selecting best parameter setting.
- The final parameters for each model/algorithm. ٠
- The computing infrastructure used for running the experiment (hardware and software), such as which software and version (libraries, frameworks, operating system etc), processing units (GPU/CPU), memory and more.
- For each reported result, the number of algorithm-runs it is averaged over and its variance.

For data used in the paper, check the following:

- For closed datasets, describe the dataset.
- For a new dataset, deposit it to a public repository with a description and metadata.
- For a new dataset, release it with a license that allows free usage for research purposes.
- All open datasets are cited.

For all code, check the following:

- All source code required for conducting the experiment is sha
- The version of the code used for conducting the experiments
- A license is added with the source code to allow free usage for

For the paper, check the following:

- Claims being investigated are stated clearly.
- For theoretical papers, complete proofs are provided (for example in the appendix). •
- Assumptions and limitations are identified.
- A conceptual outline and pseudo code describing the AI method is given.
- Statements about how the results substantiate the claims.

Source: https://folk.idi.ntnu.no/odderik/reproducibility_guidelines.pdf

ared and cited.
is specified.
or research purposes.

Recommendations	Descriptions of experiments in a publication should:
14.	Explicitly present the hypotheses to be assessed, before other details concerning the empirical study are present
15.	Present the predicted outcome of the experiment, ba on beliefs about the AI method and its application
16.	Include the experiment design (parameters and conditions to be tested) and its motivation, such as wh specific number of tests or data points are used based the desired statistical significance of results and availability of data
17.	Identify and describe the measure and metrics
18.	Provide the evaluation protocol
19.	Share the results
20.	Describe the results and the analysis
21.	Be described as a workflow that summarizes how experiment is executed and configured
22.	Include documentation on workflow executions execution traces that provide parameter settings a initial, intermediate, and final data
23.	Specify the hardware used to run the experiments
24	Be cited and published separately when complex, so t others can unequivocally refer to the individual portion of the method that they reuse or extend

Recommendations	Data mentioned in a publication should:
1.	Be available in a shared community repository, so anyone can access it
2.	Include basic metadata, so others can search and understand its contents
3.	Have a license, so anyone can understand the conditions for reuse of the d
4.	Have an associated digital object identifier (DOI) or persistent URL (PURL) data is available permanently
5.	Be cited properly in the prose and listed accurately among the references can identify the datasets unequivocally and data creators can receive creative work

Recommendations	Source code used for implementing an AI method and executing an experiment should:
6.	Be available in a shared community repository, so anyone can access it
7.	Include basic metadata, so others can search and understand its contents
8.	Include a license, so anyone can understand the conditions for use and extension software
9.	Have an associated digital object identifier (DOI) or persistent URL (PURL) for t used in the associated publication so that the source code is permanently available
10.	Be cited and referenced properly in the publication so that readers can identify to unequivocally and its creators can receive credit for their work

Recommendations	AI methods used in a publication should be:
11.	Presented in the context of a problem description that clearly identifies what problem they are intended to solve
12.	Outlined conceptually so that anyone can understand their foundational conceptually
13	Described in pseudocode so that others can understand the details of how th

(Gundersen, Gil and Aha, Al Magazine, 2018)

ncepts ey work

AAAI Reproducibility Checklist

Four sections:

- 1. The paper
- 2. Theoretical contributions
- 3. Data sets
- 4. Computational experiments

Source: https://aaai.org/Conferences/AAAI-22/reproducibility-checklist/

36th AAAI Conference on Artificial Intelligence Vancouver, BC, Canada February 22 - March 1, 2022

PROGRAM ~ CALLS ~ ORGANIZATION `

Reproducibility Checklist

Unless specified otherwise, please answer "yes" to each question if the relevant information is described either in the paper itself or in a technical appendix with an explicit reference from the main paper. If you wish to explain an answer further, please do so in a section titled "Reproducibility Checklist" at the end of the technical appendix.

his paper

AAAI-22

- clearly states what claims are being investigated (yes/partial/no)
- explains how the results substantiate the claims (yes/partial/no)
- explicitly identifies limitations or technical assumptions (yes/partial/no)
- includes a conceptual outline and/or pseudocode description of AI methods introduced (yes/partial/no/NA)

Does this paper make theoretical contributions? (yes/no)

yes, please complete the list below.

- All assumptions and restrictions are stated clearly and formally. (yes/partial/no)
- All novel claims are stated formally (e.g., in theorem statements). (yes/partial/no)
- Proofs of all novel claims are included. (yes/partial/no)
- Proof sketches or intuitions are given for complex and/or novel results. (yes/partial/no)
- Appropriate citations to theoretical tools used are given. (yes/partial/no)

Does this paper rely on one or more data sets? (yes/no)

If yes, please complete the list below.

- All novel datasets introduced in this paper are included in a data appendix (yes/partial/no/NA)
- All novel datasets introduced in this paper will be made publicly available upon publication of the paper with a license that allows free usage for research purposes (yes/partial/no/NA)
- All datasets drawn from the existing literature (potentially including authors' own previously published work) are accompanied by appropriate citations (yes/no/NA)
- All datasets drawn from the existing literature (potentially including authors' own previously published work) are publicly available (yes/partial/no/NA)
- All datasets that are not publicly available are described in detail (yes/partial/no/NA)

Does this paper include computational experiments? (yes/no)

If yes, please complete the list below.

- All source code required for conducting experiments is included in a code appendix (yes/partial/no).
- All source code required for conducting experiments will be made publicly available upon publication of the paper with a license that allows free usage for research purposes (yes/partial/no)
- If an algorithm depends on randomness, then the method used for setting seeds is described in a way sufficient to allow replication of results. (yes/partial/no/NA)
- This paper specifies the computing infrastructure used for running experiments (hardware and software), including GPU/CPU models; amount of memory; operating system; names and versions of relevant software libraries and frameworks. (yes/partial/no)
- This paper formally describes evaluation metrics used and explains the motivation for choosing these metrics (yes/partial/no)
- This paper states the number of algorithm runs used to compute each reported result (yes/no)
- Analysis of experiments goes beyond single-dimensional summaries of performance (e.g., average; median) to include measures of variation, confidence, or other distributional information (yes/no)
- This paper lists all final (hyper-)parameters used for each model/algorithm in the paper's experiments (yes/partial/no/NA)
- This paper states the number and range of values tried per (hyper-)parameter during development of the paper, along with the criterion used for selecting the final parameter setting. (yes/partial/no/NA)

The paper

- Claims are clearly stated.
- substantiate the claims. and or technical assumptions.
- Explain how the results Explicitly identify limitations
- Include conceptual outline/pseudocode of Al methods introduced.

NTNU

Research Protocol

Project Description

Short, high-level, overview

Research Questions

What is being investigated? What are the main research questions you are asking?

Why is the problem important, has anyone else said so? Briefly review previous research on each research question.

What is your contribution? How is your research topic different from what has been done before?

Methodology

What do you intend doing? Briefly describe the methods that you will use to answer your research questions.

Why is this strategy being adopted? Why is this necessary for your study?

Work Detail

Decide on the stages of the project and the dependencies between them. Compile a project plan.

- Risks (e.g., delays in obtaining key resources) and Risk Management Strategies.
- Timeline, including Gantt chart. Use specific dates so that you finish on time.
- Resources required (equipment, people, special software etc)
- Deliverables
- Milestones (which should refer to the Timeline)

Evaluation of Research Questions

You should have a plan for testing your system when it is complete. Work this out now; everything will be wasted if you finish your implementation but cannot evaluate your "advance" convincingly. Indicate the interpretation and conclusions that you will place upon the results. What difference will they make? Indicate the implications of your research for current theory and practice.

Anticipated Outcomes

- What might we expect the outcomes of your project to be? What do you expect to find? The aim here is not to anticipate your study but rather to give an outline of what you envisage.
- Major software artefacts to be produced, their key features and major design challenges.
- Key success factors how will you judge whether the project has succeeded or not?

Bibliography and Previous Systems

List the main sources on which your research will be based. In the proposal we want a preliminary outline of the key works. All work must be properly cited.

As your work progresses you have to show that you have read the relevant papers and books and understand the field. You should show that you know which important contributions are and how they are related and may be grouped. You should know where the concepts you use were first described.

Research Questions Clearly state what you are investigating?

Methodology How do you conduxt your investigation? Describe your experiment.

Evaluation of Research Questions What is the best way to evaluate the outcome of the experiment? Explain it.

Anticipate Outcomes Make a prediction. What do you expect and why?

Source: https://www.cs.uct.ac.za/teaching/forms/researchProposalGuide 2007.pdf

How well is data documented?

- We know we should not train and test on the same data.
- Is Outcome Reproducible an option if we do not know which samples were used for what?
- Can only check if Outcome \bullet Reproducible if results are shared.

The order a machine learning algorithm is fed training samples can affect the performance.

An Unbiased Look at Dataset Bias

Selection bias *Does the dataset represent a fair sampling of the world?*

Capture Bias *Are the samples represented fairly (centered object, handle direction of mugs?)*

Negative bias *Does the data set contain negative examples as well?*

(Torralba and Efros 2011

Other issues

- Data version: \bullet
 - Are there different versions of the same dataset?
 - GluonTS.

 - number of samples).
- Large dataset: \bullet
 - _____ not possible.
- Concept drift: \bullet
 - The real changes and datasets are static. _____
 - What was true one day is not true the next.

Some software libraries provide standard datasets as well i.e. seaborn and

Sometimes these differ from the original ones. Cite the correct version. Sometimes the reported data is not the same as the published data (different

Webscale datasets might not be stored after analysis. Outcome reproducibility

If the dataset is not shared it is impossible to know whether any differences are caused by concept drift or other issues related to the quality of the research.

Which Conclusions Can Be Drawn?

selection

Population

Treatment

No treatment

selection

Sample

Hardware and Ancillary Software

TABLE 1. Computing environment including FORTRAN compilers, parallel communication libraries, and optimization levels of the compiler. Identical results are marked by a symbol. Ten ensemble members with different software system are highlighted in boldface.

Name	Machine	FORTRAN compiler	Parallel communication library	Optimization level	Mark
EXP1	KISTI SUN2	INTEL 11.1	openmpi 1.4	03	
	KISTI SUN2	INTEL 11.1	mvapich2 1.5	03	
EXP2	KISTI SUN2	INTEL 11.1	mvapich1 1.2	03	0
	KISTI SUN2	INTEL 11.1	openmpi 1.4	04	
EXP3	KISTI SUN2	INTEL 11.1	openmpi 1.4	02	\triangle
EXP4	KISTI SUN2	INTEL 11.1	openmpi 1.4	01	\triangleleft
EXP5	KISTI SUN2	INTEL 11.1	openmpi 1.4	00	
EXP6	KISTI SUN2	PGI 9.0.4	openmpi 1.4	O2 (-fastsse)	
	KISTI SUN2	PGI 9.0.4	mvapich2 1.5	O2 (-fastsse)	
	KISTI SUN2	PGI 9.0.4	mvapich1 1.2	O2 (-fastsse)	
	KISTI SUN2	PGI 8.0.6	mvapich1 1.2	O2 (-fastsse)	
	YSU Cluster	PGI 10.6	mvapich1 1.2	O2 (-fastsse)	
	YSU Cluster	PGI 10.6	mvapich1 1.2	O3 (-fastsse)	
EXP7	YSU Cluster	PGI 10.6	mvapich1 1.2	01	•
EXP8	YSU Cluster	PGI 7.1.6	mvapich1 1.2	O2 (-fastsse)	
EXP9	KISTI IBM 1	XLF 10.1	_	03	*
	KISTI IBM 2	XLF 12.1		03	*
	KISTI IBM 1	XLF 10.1		04	*
EXP10	KISTI IBM 1	XLF 10.1	_	02	٠
	KISTI IBM 1	XLF 10.1	_	01	٠

Code Version

Ξ

What data does R keep if you run *distinct(data, Subject)*? Menu Depends. When did you last update {dplyr}? condition mediator Subject [95] Groundhog: Addressing The Threat 70 treatmen Before June 24 2016 1154 555 3 control condition mediator Subject dv 47888 110 placebo treatment 11543 70 5 That R Poses To Reproducible Research treatment 70 11543 5 555 control 555 6 3 control After June 24 2016 110 3 47888 placebo Subject 47888 110 3 placebo 11543 Posted on January 5, 2021 by Uri Simonsohn 555 47888

R, the free and open source program for statistical computing, poses a substantial threat to the reproducibility of published research. This post explains the problem and introduces a solution.

The Problem: Packages

R itself has some reproducibility problems (see example in this footnote [1]), but the big problem is its packages: the addon scripts that users install to enable R to do things like run meta-analyses, scrape the web, cluster standard errors, format numbers, etc. The problem is that packages are constantly being updated, and sometimes those updates are not backwards compatible. This means that the R code that you write and run today may no longer work in the (near or far) future because one of the packages your code relies on has been updated. But worse, R packages depend on other packages. Your code could break after a package you don't know you are using updates a function you have never even used.

When you run the code later, you might get different results!

Computational experiments II

- This paper formally describes **evaluation metrics** used and explains the motivation for choosing these metrics.
- This paper states the number of algorithm runs used to compute each reported result.
- Analysis of experiments goes beyond single-dimensional summaries of performance (e.g., average; median) to include measures of variation, confidence, or other distributional information
- This paper lists all final (hyper-)parameters used for each model/algorithm in the paper's experiments.
- This paper states the number and range of values tried per (hyper-)parameter during development of the paper, along with

Deep Learning that Matters

- Hyperparameter search will have a huge effect on results. Ranges rarely documented properly.
- Simple changes in network architecture can have make large changes to result.

Different implementations of same baseline algorithm can yield very different results.

Algorithm Runs and Variation I

Ran the same experiment 100 times. Only difference was which seeds we used to initialize the pseudorandom number generator

Algorithm Runs and Variation II

KDE used to smooth out the variance of a selection of seeds. See how different the average MAPE scores for those seeds will be. Assuming a similar distribution for our baseline, we can manipulate results by selecting the best set of 5 seeds for our algorithm and the 5 worst seeds for our baseline.

Experiment: MNIST Classification I

- Same experiment conducted 20 times on four different machine learning platforms.
- Code = same
- Data = same
- HW = **!same**
- Ancillary SW = *!same*

Experiment: MNIST Classification II

When models are wrong, how many are wrong?

CPU, random seed *not* fixed

CPU, random seed fixed

Experiment: MNIST Classification III

CPU, random seed *not* fixed

When models are wrong, how many different classes do they see?

CPU, random seed fixed

(Gundersen, Shamsaliei and Isdahl, forthcoming)

THE VALUE AND CHALLENGES OF TRANSPARENT RESEARCH PART IV

The Ten Years Reproducibility Challenge

TEN VERS REPRODUCIBILITY CHALLENGE RESCIENCE SPECIAL ISSUE FREE TO READ - FREE TO PUBLISH

Would you dare to run the code from your past self?

(the one that does not answer mail)

S U B M I S S I O N D E A D L I N E 0 1 / 0 4 / 2 0 2 0 **http://rescience.github.io/ten-years** In association with Inria, CNRS, Software Heritage, ReScience, Comité pour la Science Ouverte, URFIST Bordeaux & Mission de la pédagogie et du numérique pour l'enseignement supérieur. **Contact:** nicolas.rougier@inria.fr "Programming languages evolve, as do the computing environments in which they run, and code that works flawlessly one day can fail the next."

Nicolas Rougier, Nature, 2020

PoV of Original Researchers

Incereased documentation efforts

(Gundersen, Gil and Aha, 2018)

PoV of Independent Researchers

Incereased trust in the original study's results

(Gundersen, Gil and Aha, 2018)

Reproducibility Experiment

Success: 20%

Partial success: 13%

Failure: 23%

No result: 17%

Filtered out (R3): 27%

(Gundersen et al, forthcoming)

The value of sharing both code and data

We tried to reproduce 30 of the topcited papers from 2012, 2014 and 2016. These are the results: *Sharing* both code and data is really effectfull.

83%

Success (green), Partial success (orange), Failure (red) and no result (grey) when reproducing experiments with and without code. Each box represents an aggregate of the experiments reported in one paper (most cited Al papers from Scopus).

(Gundersen et al, forthcoming)

CONCLUSION: WHAT IF YOU CANNOT DO EVERYTHING? PART V

Important to Remember

State of the Art: Reproducibility in Artificial Intelligence

Odd Erik Gundersen and Sigbjørn Kjensmo

Department of Computer Science Norwegian University of Science and Technology

Hacker News new | threads | past | comments | ask | show | jobs | submit

- ▲ State of the Art: Reproducibility in Artificial Intelligence [pdf] (aaai.org) 43 points by capablemonkey on Oct 6, 2018 | hide | past | favorite | 6 comments
- Bac

ing

duci

met

Hyp

to re

have

and

proc

and

beer

fron

veye

men

and

of th

whil

▲ sgt101 on Oct 6, 2018 [-]

I think that the result is overcooked. Their hypothesis 1 is somewhat falsifiable in that I don't think that there is a widespread reproducibility crisis. I have been unable to reproduce results a couple of times in my career, but I think that each time that was due to naughtiness (deliberate) on the authors part or incompetence by me. Almost always you can reproduce and when I have run into trouble I've found that the authors almost always help out (most people are just delighted that you are interested!) On the other hand this paper is very useful in that I think it will be used to establish better criteria for papers in the future. I often reject papers because they make no claim and have no results, contribution or conclusions (this makes reviewing them quick so I really like papers like this !) I think that it would be harsh to outright reject a paper because the hardware set up is poorly documented, but it would be reasonable to ask for that change before publication (for example). I agree with the authors that their criteria are useful.

One issue though, open sourcing software is a good aspiration, but it's not always possible due to IP and licensing - also export controls in some cases (not always US -> other places too). If the community insists on opensource pre-publication some important stuff is not going to get published.

producibility scores decrease with increased documentation requirements. Improvement over time is found. Conclusion: Both hypotheses are supported.

same result as the original researchers, then they refute the hypothesis" (Oates 2006, p. 285). Hence, the inability to reproduce results affects the trustworthiness of science. To ensure high trustworthiness of AI and machine learning re-

Comment from Hacker News

andve *et al.* 2013; d focus on reproption of data and et al. 2013). Still, ucibility see little time required to lutions (Gent and gues that automaor machine learna computer. Deat is reproducible ficial intelligence kkens *et al.* 2013;

producibility; "if ment and get the

Many people believes that the reason that they are not able to reproduce results is their own incompetency.

This leads to false claims not being refuted!

Newton did not share code and data Writing a good paper that describes the experiment well and is fully transparent is most important!

- \bullet
- Sharing is more important than good documentation.

If you do not have time to document and tidy up the code and data, it is still better to share the code and data than not to.

SUCCESSFUL RESEARCH IN AI

Standing on the Feet of Giants Odd Erik Gundersen, dr. philos.

Chief AI Officer, TrønderEnergi AS Adjunct Associate Professor, NTNU odderik@ntnu.no

Norwegian Open AI Lab

References

• Tian, Y., Ma, J., Gong, Q., Sengupta, S., Chen, Z., Pinkerton, J., & Zitnick, C. L. (2019). Elf opengo: An analysis and open reimplementation of alphazero. ICML. arXiv preprint arXiv:1902.04522.

Research

- 2018
- **On Reproducible AI** O. E. Gundersen, Y. Gil and D. W. Aha, AI Magazine, Fall 2018. \bullet
- **Standing on the Feet of Giants** O. E. Gundersen, Al Magazine, Winter 2019.
- Isdahl and O. E. Gundersen, eScience 2019.
- The Reproducibility Crisis Is Real, O. E. Gundersen, Al Magazine, 41(3), 103-106, 2020.
- **The Case Against Registered Reports**, O. E. Gundersen, Al Magazine, Spring 2021.
- What We Learned When Reproducing the Most Cited Al Research, O. E. Gundersen, O. Cappelen, N. Grimstad, M. Mølnå, forthcoming.

Odd Erik Gundersen <u>odderik@ntnu.no</u>

State of the Art: Reproducibility in Artificial Intelligence O. E. Gundersen and S. Kjensmo, AAAI

A Method for Assessment and a Survey of the Reproducibility Support of ML Platforms, R.

References

- https://arxiv.org/abs/1902.04522
- P. Henderson et al. (2018). Deep Reinforceement Learning that Matters, AAAI 2018.
- M. Baker (2016), Is There a Reproducibility Crisis?, Nature, 2016.
- Monthly Weather Review, 141(11):4165–4172, 2013.
- IEEE conference on computer vision and pattern recognition (pp. 770-778).

- Surveys, 1991.

Odd Erik Gundersen <u>odderik@ntnu.no</u>

Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., ... & Hassabis, D. (2016). Mastering the game of Go with deep neural networks and tree search. Nature, 529(7587), 484-489. Tian et al (2019), ELF OpenGo: An Analysis and Open Reimplementation of AlphaZero, ICML 2019, URL:

Song-You Hong, Myung-Seo Koo, Jihyeon Jang, Jung- Eun Esther Kim, Hoon Park, Min-Su Joh, Ji-Hoon Kang, and Tae-Jin Oh (2013), An evaluation of the software system dependency of a global atmospheric model.

He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In Proceedings of the

Pineau, J., Vincent-Lamarre, P., Sinha, K., Larivière, V., Beygelzimer, A., d'Alché-Buc, F., ... & Larochelle, H. (2020). Improving Reproducibility in Machine Learning Research (A Report from the NeurIPS 2019 Reproducibility Program). *arXiv preprint arXiv:2003.12206*.

Raff, E. (2020). Research Reproducibility as a Survival Analysis. arXiv preprint arXiv:2012.09932.

Torralba, A., & Efros, A. A. (2011). Unbiased look at dataset bias. In CVPR 2011 (pp. 1521-1528). IEEE.

D. Goldberg (1991). What Every Computer Scientist Should Know About Floating-Point Arithmetic, Computing

Resources

- **Ten Year Reproducibility Challenge**, Rescience C,URL: https://rescience.github.io/ten-years/
- \bullet
- https://www.cs.mcgill.ca/~jpineau/ReproducibilityChecklist. Pdf

Challenge to scientists: does your ten-year-old code still run?, J. M. Perkel, Nature (2020), URL: https://www.nature.com/articles/d41586-020-02462-7

General Reproducibility Guidelines for AI, Gundersen, Gil, Mausam, 2020, URL: https://folk.idi.ntnu.no/odderik/reproducibility guidelines.pdf

The Machine Learning Reproducibility Checklist, J. Pineau, 2020, URL:

AlphaGo – the Movie, Youtube, URL: https://youtu.be/WXuK6gekU1Y