

# What I talk about when I talk about when I talk about reproducibility

Odd Erik Gundersen, dr. philos. Chief Al Officer, TrønderEnergi AS Adjunct Associate Professor, NTNU odderik@ntnu.no



NTNU Norwegian University of Science and Technology





## INTRODUCTION





#### 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)







#### Computer Science





### ICLR 2018 Reproducibility Challenge

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





(J. Pineau, ICLR keynote, 2018)





### Repeatability in Comp. Sys. Research

Summary of the study's results. Blue numbers represent papers we excluded from the study, green numbers papers we determined to be weakly repeatable, red numbers papers we determined to be non-repeatable, and orange numbers represent papers for which we could not conclusively determine repeatability (due to our restriction of sending at most one email request per author).



- Analyzed 601 ACM papers.
- Out of these
- Locate and build source code.
- Able to for 32.3% w.o. communicating with authors.
- Increase to 48.3% with communication.





### **Journal Policy Effectiveness Analysis**

Table 1. Responses to emailed req	<ul> <li>Science policy to</li> </ul>		
Type of response	Count	Percent, %	
Did not share data or code:			data.
Contact another person	20	11	Doguoatad fram 90
Asked for reasons	20	11	• Requested from ZU
Refusal to share	12	7	nanare from
Directed back to supplement	6	3	μαμεις ποιπ
Unfulfilled promise to follow up	5	3	Science
Impossible to share	3	2	
Shared data and code	65	36	Obtained artifacte
Email bounced	3	2	· Uplained alliadis
No response	46	26	from 44%.
			<ul> <li>Able to reproduce</li> </ul>

26%.



### PART 1 UNDERSTANDING REPRODUCIBILITY



### Reproducibility SO WHAT IS IT?







#### The Scientific Method in Empirical AI Research









### **Types of studies**

#### Hypothesis generating

#### Observatory

Manipulatory

#### **Assessment studies**

#### **Exploratory studies**

#### 



(P. R. Cohen, MIT Press, 1995)



#### Which conclusions can we draw?



selection

#### Population

#### Treatment

#### No treatment

selection

Sample



### Defining Reproducibility I



#### **Reproducibility Spectrum**

Code and data Linked and executable code and data Full replication

Gold standard

(R. D. Peng, Science, 2011)





# **Defining Reproducibility II**

**Replication** is to re-run the experiment with code and data provided by the author.

**Reproduction** implies both replication and the regeneration of findings with at least some independence from the [original] code and/or data.





# Defining Reproducibility III

the same data and tools, to obtain the same results.

reanalysis of the original study.

- **Methods reproducibility:** The ability to implement, as exactly as possible, the experimental and computational procedures, with
- **Results reproducibility:** The production of corroborating results in a new study, having used the same experimental methods.
- **Inferential reproducibility:** The drawing of qualitatively similar conclusions from either an independent replication of a study or a





### **Definition of Reproducibility**

Reproducibility in empi an **independent** researed results using the same documentation made



- Reproducibility in empirical AI research is the ability of
- an independent research team to produce the same
- results using the same AI method based on the
- documentation made by the original research team.

### In order to reproduce results WHAT MUST BE DOCUMENTED?







### **Reproducibility crisis**



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



#### The Scientific Method in Empirical AI Research







### Documentation

- Method (text): Description of AI method
- Used for testing hypotheses.

(system/algorithm), study design, experiment description - human to human, abstract concepts.

Data: Represents the world the AI method operates in.

Experiment (software): Al method code + experiment code + experiment setup + HW + SW + analysis code



### Degree of Reproducibility

	Method	Data	Experiment
R1			
R2			
R3			





Factor	Variable	
	Problem	
	Objective	
	Research method	
Method	Research questions	
	Pseudocode	
	Hypothesis	
	Prediction	
	Experiment setup	
	Training data	
	Validation data	
Data	Test data	
	Results	
Experiment	Method source code	
	Experiment source code	
	Software dependencies	
	Hardware	



Description
Is there an explicit mention of the problem the research seeks to solve?
Is the research objective explicitly mentioned?
Is there an explicit mention of the research method used (empirical, theoretical)?
Is there an explicit mention of the research question(s) addressed?
Is the AI method described using pseudocode?
Is there an explicit mention of the hypotheses being investigated
Is there an explicit mention of predictions related to the hypotheses?
Are the variable settings shared, such as hyperparameters?
Is the training set shared?
Is the validation set shared?
Is the test set shared?
Are the relevant intermediate and final results output by the Al program shared?
Is the AI system code available open source?
Is the experiment code available open source?
Are software dependencies specified?
Is the hardware used for conducting the experiment specified?





### **QUANTIFYING REPRODUCIBILITY**





### **Quantifying Reproducibility**

## $R1D(e) = \frac{\delta_1 Method}{}$

 $R2D(e) = \frac{\delta_1 l}{d}$ 

R3D(

$$\frac{d(e) + \delta_2 Data(e) + \delta_3 Exp(e)}{\delta_1 + \delta_2 + \delta_3}$$
$$\frac{Method(e) + \delta_2 Data(e)}{\delta_1 + \delta_2}$$
$$e) = Method(e)$$



### **A Normalized Metric**



(Gundersen, Kjensmo, AAAI, 2018)





## WHAT WE GAIN





#### We Can Specify How Well Research is **Documented**

#### Method



#### Data

#### Experiment





### We Can Measure Improvement







### We Can Compare Research: Papers

Id	Title		Y ear	Hours
10				spent
1	Measuring the Objectness of Image Windows [26]	$\mathbf{R1}$	2012	40
2	Generalized Correntropy for Robust Adaptive Filtering [27]	R2-D	2016	40
2	Development and investigation of efficient artificial bee colony algorithm	ם פס	2012	40
5	for numerical function optimization [28]		2012	40
4	Blind Image Quality Assessment: A Natural Scene Statistics Approach	D1	2012	25
4	in the DCT Domain $[29]$	RΙ	2012	20
5	Cooperatively Coevolving Particle Swarms for Large Scale Optimization		2012	40
9	[30]	<u>π</u> 2-D	2012	40
6	Learning Sparse Representations for Human Action Recognition [31]	R2-D	2012	40
7	Visualizing and Understanding Convolutional Networks [32]	R2-D	2014	40
	iSuc-PseOpt: Identifying lysine succinylation sites in proteins by incor-			
8	porating sequence-coupling effects into pseudo components and optimiz-	R2-D	2016	22
	ing imbalanced training dataset [33]			
0	A modified Artificial Bee Colony algorithm for real-parameter optimiza-	ם פס	2012	40
9	tion $[34]$	<u>π</u> 2-D	2012	40
10	RASL: Robust alignment by sparse and low-rank decomposition for lin-	D1	2012	10
	early correlated images $[35]$	ΠI	2012	10





### We Can Compare Research: Conferences

Conference	$R1D \pm \varepsilon$	$R2D \pm \varepsilon$	$R3D \pm \varepsilon$
IJCAI 2013	$0.20 \pm 0.02$	$0.20 \pm 0.03$	$0.24 \pm 0.04$
AAAI 2014	$0.21 \pm 0.02$	$0.26 \pm 0.03$	$0.28 \pm 0.04$
IJCAI 2016	$0.30\pm0.03$	$0.30\pm0.04$	$0.29 \pm 0.04$
AAAI 2016	$0.23 \pm 0.02$	$0.25\pm0.04$	$0.24 \pm 0.04$
Total	$0.24 \pm 0.01$	$0.25 \pm 0.02$	$0.26 \pm 0.02$





### We Can Compare Research: Groups

#### **Academia versus Industry**

Test



#### Method



Data

#### Experiment

(Gundersen, AI Magazine, forthcoming)





### We Can Compare Software Frameworks



(Isdahl and Gundersen, forthcoming)





#### We Could Empirically Find What Entails Well-**Documented Research**

Factor	Variable	Description
Method	Problem	Is there an explicit mention of the problem the research seeks to solve?
Objective		Is the research objective explicitly mentioned?
	Research method	Is there an explicit mention of the research method used (empirical, theoretical)?
	Research questions	Is all explicit mention of the research question(s) dressed?
	Pseudocode	the AI methescribed using pseudocode?
Data	Training data	ne training shared?
	Validation data	Is the validat set shared?
	Test data	Is the test sanared?
	Results	Are the revant intermediate and final results output by the AI program shared?
Experiment	Hypothesis	Is there in explicit mention of the hypotheses being investigated?
	Prediction	Is there an explicit mention of predictions related to the hypotheses?
	Method source code	Is the stem code available open source?
	Hardware	Is the hardware used for conducting the experiment specified?
	Software dependencies	Are software dependencies specified?
	Experiment setup	Are the variable settings shared, such as hyperparameters?
	Experiment source code	Is the experiment code available open source?



### **Compute the Likelihood of Success?**





# We Can Set the Bar Based on What We Want to Achieve



500motivators.com

### CASE STUDY HOW WELL IS AI RESEARCH DOCUMENTED?




# **Experiment**

• We surveyed 400 papers.

AAAI 2016, IJCAI 2013 and IJCAI 2016.

the reproducibility.



# 100 papers from each installment of AAAI 2014,

#### Six reproducibility metrics proposed for quantifying

(Gundersen, Kjensmo, AAAI, 2018)





# Degree of Reproducibility

	Method	Data	Experiment
R1			
R2			
R3			





### **Results I: Factors and Variables**

#### Method



Data

#### Experiment





# **Results II: Reproducibility Degree**



(Gundersen, Kjensmo, AAAI, 2018)





# Results III: Change over Time







## **Results IV: Industry vs Academia**



#### Method



Data

Test

#### Experiment

(Gundersen, AI Magazine, forthcoming)





## **Results V: Industry vs Academia**









#### PART 2 CAUSES OF IRREPRODUCIBILITY



#### 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						
			10			100-
	0	20	40	60	80	100%





#### MISSING DATA

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



**46** 





#### **Deep Reinforcement Learning that Matters**





- Non-determinism in standard benchmark environments and
- Variance intrinsic to the method
- Cause irreproducible results





### **Deterministic Implementations for Reproducibility in DRL**





(f) Minibatch

- Non-determinism in training process.
- Deterministic implementation of Qlearning.
- Measure impact of different sources of nondeterminism.
- Different sources have huge impact on performance.





# Are GANs created equal?



Measure = Recall DRAGAN MM GAN NS GAN NGAN WGAN G

 $\bullet$ 

- Study on models and evaluation measures.
  - Most models can reach same performance given hyperparameter optimization and random restarts.
  - Suggests more systematic and objective evaluation procedures.





#### Software Dependency of Weather Model

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	<b>KISTI SUN2</b>	<b>INTEL 11.1</b>	openmpi 1.4	03	
	<b>KISTI SUN2</b>	<b>INTEL 11.1</b>	mvapich2 1.5	O3	
EXP2	<b>KISTI SUN2</b>	<b>INTEL 11.1</b>	mvapich1 1.2	03	0
	<b>KISTI SUN2</b>	<b>INTEL 11.1</b>	openmpi 1.4	O4	
EXP3	<b>KISTI SUN2</b>	<b>INTEL 11.1</b>	openmpi 1.4	02	$\triangle$
EXP4	<b>KISTI SUN2</b>	<b>INTEL 11.1</b>	openmpi 1.4	01	$\triangleleft$
EXP5	<b>KISTI SUN2</b>	<b>INTEL 11.1</b>	openmpi 1.4	00	
EXP6	<b>KISTI SUN2</b>	PGI 9.0.4	openmpi 1.4	O2 (-fastsse)	
	<b>KISTI SUN2</b>	PGI 9.0.4	mvapich2 1.5	O2 (-fastsse)	
	<b>KISTI SUN2</b>	PGI 9.0.4	mvapich1 1.2	O2 (-fastsse)	
	<b>KISTI SUN2</b>	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	*
	<b>KISTI IBM 2</b>	XLF 12.1		O3	*
	<b>KISTI IBM 1</b>	XLF 10.1		O4	*
EXP10	<b>KISTI IBM 1</b>	XLF 10.1		02	٠
	KISTI IBM 1	XLF 10.1		01	٠



# REPRODUCING THE MOST CITED AI RESEARCH

Case study





# Experiment

• We selected 30 papers to reproduce

based on numbers from Scopus.

Structured research procedure.



# Ten most cited AI papers from 2012, 2014 and 2016

(Gundersen et al, forthcoming)





## **Research Procedure**

- Reproduce research classified as R1 and R2 reproducible.
- Time-boxed the work put into each research paper to 40 hours effective work time.
- Stopping criteria (computing resources, paywall data sets, only qualitative results presented).





# Degree of Reproducibility

	Method	Data	Experiment
R1			
R2			
R3			





# **Results: Reproducibility Degree**

8



(Gundersen et al, forthcoming)





## Results: Outcome per paper





Success: 3%

Partial success: 30%

Failure: 30%

No result: 23%

Filtered out (R3): 27%

(Gundersen et al, forthcoming)





# **Top Six Causes of Failure**

- Aspect of *implementation* not described or ambiguous (R2).
- Aspect of *experiment* not described or ambiguous (R2). Not all *hyper-parameters* are specified (R2).  $\bullet$
- Mismatch between data in paper and available online (R1+R2).
- Method code shared, experiment code not shared (R1). Method not described with enough detail (R2). •





#### The Scientific Method in Empirical AI Research







#### Part 3 RECOMMENDATIONS



### **Recommendations in Al Magazine**

- Author checklist of 24 practical recommendations lacksquare
- Four groups:  $\bullet$ 
  - Data
  - Source code \_\_\_\_\_
  - Al Methods
  - Experiments \_\_\_\_\_
- Summary:
  - Open Science share data, code and procedures.
  - Build digital scholarship
  - Version code and data!

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



#### Data

Recommendations	Data mentioned in a pi
1.	Be available in a share
2.	Include basic metadat
3.	Have a license, so any
4.	Have an associated dig data is available perm
5.	Be cited properly in t can identify the datas work

#### ublication should:

- ed community repository, so anyone can access it
- ta, so others can search and understand its contents
- one can understand the conditions for reuse of the data
- gital object identifier (DOI) or persistent URL (PURL) so that the anently
- the prose and listed accurately among the references, so readers sets unequivocally and data creators can receive credit for their





## Source code

Recommendations	Source code used for implem
6.	Be available in a shared co
7.	Include basic metadata, so
8.	Include a license, so anyo software
9.	Have an associated digita used in the associated pub
10.	Be cited and referenced pr unequivocally and its crea

enting an AI method and executing an experiment should:

ommunity repository, so anyone can access it

o others can search and understand its contents

one can understand the conditions for use and extension of the

l object identifier (DOI) or persistent URL (PURL) for the version plication so that the source code is permanently available

roperly in the publication so that readers can identify the version ators can receive credit for their work



## Al Methods

Recommendations	AI methods used in a publ
11.	Presented in the contempoblem they are intended
12.	Outlined conceptually
13	Described in pseudocod

lication should be:

xt of a problem description that clearly identifies what led to solve

so that anyone can understand their foundational concepts

de so that others can understand the details of how they work





#### Experiments

-	Recommendations	Descripti
	14.	Explicitl other de
	15.	Present on belie
	16.	Include conditio specific the des availabil
	17.	Identify
	18.	Provide
	19.	Share th
	20.	Describe
	21.	Be desc experim
	22.	Include executio initial, in
	23.	Specify t
	24	Be cited others c of the m

#### tions of experiments in a publication should:

etails concerning the empirical study are presented

the predicted outcome of the experiment, based efs about the AI method and its application

the experiment design (parameters and the ons to be tested) and its motivation, such as why a number of tests or data points are used based on sired statistical significance of results and the ility of data

v and describe the measure and metrics

the evaluation protocol

he results

e the results and the analysis

cribed as a workflow that summarizes how the nent is executed and configured

e documentation on workflow executions or on traces that provide parameter settings and intermediate, and final data

the hardware used to run the experiments

and published separately when complex, so that can unequivocally refer to the individual portions nethod that they reuse or extend



# The ML Reproducibility Checklist

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

- external libraries.

For any theoretical claim, check if you include:

- A statement of the result.
- A clear explanation of any assumptions.
- A complete proof of the claim.

A clear description of the mathematical setting, algorithm, and/or model.

An analysis of the complexity (time, space, sample size) of any algorithm.

A link to a downloadable source code, with specification of all dependencies, including



# The ML Reproducibility Checklist

For all figures and tables that present empirical results, check if you include:

- A complete description of the data collection process, including sample size.
- A link to a downloadable version of the dataset or simulation environment.
- An explanation of any data that were excluded, description of any pre-processing step.
- An explanation of how samples were allocated for training / validation / testing.
- 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 evaluation runs.
- A description of how experiments were run.
- A clear definition of the specific measure or statistics used to report results.
- Clearly defined error bars.
- A description of results with <u>central tendency</u> (e.g. mean) & <u>variation</u> (e.g. stddev).
- A description of the computing infrastructure used.



#### Reproducibility in Computational and **Experimental Mathematics 2012 Implementation Criteria I:**

- methods.
- Necessary run parameters were given.

A precise statement of assertions to be made in the paper. Full statement (or valid summary) of experimental results. Salient details of data reduction & statistical analysis

A statement of the computational approach, and why it constitutes a rigorous test of the hypothesized assertions.

Complete statements of, or references to, every algorithm used, and salient details of auxiliary software (both research and commercial software) used in the computation.





## Reproducibility in Computational and **Experimental Mathematics 2012**

#### **Implementation Criteria II:**

- Discussion of the adequacy of parameters such as precision level and grid resolution.
- Proper citation of all code and data used, including that generated by the authors.
- Availability of computer code, input and output data, with some reasonable level of documentation.
- Avenues of exploration examined throughout development, including  $\bullet$ information about negative findings.
- Instructions for repeating computational experiments described in the  $\bullet$ article.
- Precise functions were given, with settings.  $\bullet$ Salient details of the test environment, including hardware, system software, and number of processors used.





#### Reproducibility in Computational and **Experimental Mathematics 2012 Archiving Criteria:**

- Data documented to clearly explain what each part represents.
- Data archived with significant longevity expected.
- Data location provided in the acknowledgements.
- Authors have documented use and licensing rights.
- Software documented well enough to run it and what it ought to  $\bullet$ do.
- The code is publicly available with no download requirements. There was some method to track changes/to the software, as well as some certainty that the code is securely archived.





#### PART 4 CHALLENGES





## **Open Questions for Reproducibility**

- Will it rely on documenation or result or both?
- What do we mean by the same result?
- Are the levels properly defined?
- What exactly must be documented?
- For each level?

Can we agree on a definition of reproducibility?

Does this change between reproducibility degrees?



# **PoV of Original Researchers**

#### Increased documentation efforts



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




#### **PoV of Independent Researchers**



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



## **Barriers to reproducibility**

- **Time consuming:** Proper documentation, questions from external researchers, maintenance cost.
- No incentives: Not required by publishers, grant makers, evaluating committees for research positions.
- **Risk future work:** Sharing of data, code and detailed experiment procedures will enable independent researchers to quickly build on the published research, and jeopardize possible new publications.





#### PART 5 FUTURE



## **Removing Barriers**

- Build infrastructure: Reduce the effort for individuals.
- **Provide infrastructure:** Publishers, academic institutions, and grant makers could provide the infrastructure.
- Eligibility requirements: Make reproducibility a requirement for academic positions.
- Emphasize quality: Not quantity when evaluating researchers for positions.
- Reward sharing: Not only review how many papers have been published, but also how many data sets and code repositories are shared when reviewing candidates.
- Reward reproducibility: Have others reproduced the research? Does it represent scientific knowledge?



### Scientific Paper of the Future

#### Scientific Paper of the Future

**Modern Paper** 

Text: Narrative of the method, some data is in tables, figures/plots, and the software used is mentioned

Data: Include data as supplementary materials and pointers to data repositories

#### **Reproducible Publication**

Software: For data preparation, data analysis, and visualization

**Provenance and methods:** Workflow/scripts specifying dataflow, codes, configuration files, parameter settings, and runtime dependencies

#### **Open Science**

Sharing: Deposit data and software (and provenance/workflow) in publicly shared repositories

**Open licenses:** Open source licenses for data and software (and provenance/workflow)

#### Metadata:

Structured descriptions of the characteristics of data and software (and provenance/workflow)

#### **Digital Scholarship**

**Persistent identifiers:** For data, software, and authors (and provenance/workflow)

#### **Citations:**

Citations for data and software (and provenance/workflow)

- Tutorial @ AAAI 2017
- Pratical advise on how to make your research reproducible.
- Complete tutorial available at website.





#### **Machine Learning Platforms**





#### **Machine Learning Platforms**





# In the Future, Who Are Responsible for What?

- Academic institutions: Evaluate whether open positions will be filled by scientists whos research is reproducible.
- **Publishers:** Provide infrastructure, such as code and data repositories, and guidelines on how to publish research including code and data.
- **Grant makers:** Require funded research to be reproducible and that data and code are shared.
- Scientists: Ensure that proper science is conducted.



## What if We Cannot Share?

- Many valid reasons for not sharing.
  - Privacy, cannot share private data.
  - Data set is too large.
  - Company IP.
  - Commercial software.



## **Open Questions**

- Should public money fund research for which results are not shared with the public?
- Should such a requirement apply to academia only? • What if academia collaborate with industry?
- Is it not better to publish papers that describe ideas and do not share code and data, than not publish the ideas?
- How can we ensure that industry continue to publish?
- Should we be more explicit; should we label our research with the reproducibility degree?



## SUMMARY

NTNU |Odd Erik Gundersen|odderik@ntnu.no



## **Current State**

- Most research is not reproducible.
- Research is so poorly documented that it is hard if even possible to reproduce the results.
- According to loannidis, most research findings are false.





## **Scenario I - Dark Future**

- We continue in the same vain. • Al research loses credibility, as it is do not follow
- scientific method.





## **Scenarion II - Bright Future**

- Many good tools exist that support reproducible experiments. We start using them.
- The amount of research that share data and code increases.
- We see that small changes make a huge impact and we improve further.
- Fewer dead-ends are visited and knowledge improves faster.
- Virtuous circle.



## Important Notes

- Perfectly executed research need not be reproducible!
- wrong!
- others to validate results.

• Claims made by reproducible research might be

Reproducibility is about transparency and enabling



## **Reproducibility is a core part of science**

# Are your experiments reproducible?



#### References

- Kjensmo, AAAI 2018.
- 2018.
- $\bullet$ Isdahl and O. E. Gundersen, eScience 2019.
- $\bullet$ Gundersen, O. Cappelen, M. Mølnå and N. Grimstad, forthcoming.

NTNU | Odd Erik Gundersen | odderik@ntnu.no

#### State of the Art: Reproducibility in Artificial Intelligence O. E. Gundersen and S.

**On Reproducible AI** O. E. Gundersen, Y. Gil and D. W. Aha, AI Magazine, Fall

**Standing on the Feet of Giants** O. E. Gundersen, Al Magazine, forthcoming 2019.

**Out-of-the-box Reproducibility: A survey of Machine Learning Platforms**, R.

What We Learned when Reproducing the Most Cited AI Research, O. E.



#### References

- The Machine Learning Reproducibility Checklist, J. Pineau, 2019
- Setting the Default to Reproducible, ICERM workshop report, https://icerm.brown.edu/topical\_workshops/tw12-5-rcem/icerm\_report.pdf, 2012
- An Empirical Analysis of Journal Policy Effectiveness for Computational **Reproducibility**, Stodden, Seiler and Ma, PNAS, 2018.
- An Evaluation of the Software System Dependency of a Global Atmospheric **Model**, S. HONG, M. KOO, and J. JANG, Monthly Weather Review, 2013.
- Are GANs Created Equal? A Large-scale Study Lucic, M., Kurach, K., Michalski, M., Gelly, S., & Bousquet, O., Advances in neural information processing systems. 2018.
- **1,500** scientists lift the lid on reproducibility, M. Baker, Nature News, 533(7604), 452, 2016.

#### NTNU | Odd Erik Gundersen | odderik@ntnu.no

#### **Deterministic Implementations for Reproducibility in Deep Reinforcement Learning**, P. Nagarajan, G. Warnell, P. Stone, AAAI 2019 workshop on Reproducible AI, 2019



#### References

- Learning. 2019.
- Repeatability in computer systems research, C. Collberg and T. A. Proebsting, Communications of the ACM, 2016
- Learn to Write a Scientific Paper of the Future: **Reproducible Research, Open Science, and Digital** Scholarship, Y. Gil, D. Garijo, G. Peretsman-Clement, AAAI 2017 Tutorial, http://scientificpaperofthefuture.org/
- Why Most Published Research Findings Are False, J. P. Ioannidis, *PLoS medicine*, 2(8), e124, 2005

NTNU | Odd Erik Gundersen | odderik@ntnu.no

#### **Unreproducible Research is Reproducible**, X. Bouthillier, C. Laurent, and P. Vincent, International Conference on Machine