# $R^3$ — *Responsible Replicable Research*

Margot Mieskes

University of Applied Sciences, Darmstadt

June 27, 2025

**h_da**
HOCHSCHULE DARMSTADT
UNIVERSITY OF APPLIED SCIENCES

**fbmd**
FACHBEREICH MEDIA

# Are you a researcher

# Are you a researcher

What defines a researcher?

# Are you a researcher

What defines research?

# Research

$R^3$ —
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

Community

Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
Call to Action

- creative
- systematic
- control bias
- control error

https://en.wikipedia.org/wiki/Research

# Research

$R^3$ —
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

Community

Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
Call to Action

- creative
- systematic
- control bias
- control error

replicate elements of prior projects

`https://en.wikipedia.org/wiki/Research`

# Reproducible vs. Replicable

What does Replicable/Reproducible mean?

# Definitions – kind of

$R^3$ –
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

Community

Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
Call to Action

**Replicability** or **repeatability** is a property of an *experiment*: the ability to repeat – or not – the experiment described in a study.

Cohen et al LREC 2018 https://aclanthology.org/L18-1025/

# Definitions – kind of

R³ –
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

Community

Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
Call to Action

**Reproducibility** is a property of the *outcomes* of an
experiment: arriving – or not – at the same conclusions,
findings, or values.

Cohen et al LREC 2018 https://aclanthology.org/L18-1025/

# Dimensions of Reproducibility

R³ —
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

Community

Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
Call to Action

Conclusion: Broad induction that is made based on the results of the reported research.

Finding: Relationship between the values for some reported figure of merit with respect to two or more dependent variable.

Value: A number, whether measured (e.g. acount of false positives) or calculated (e.g. a standard deviation).

Cohen et al LREC 2018 https://aclanthology.org/L18-1025/

# Research Artifacts

- Data
- Code
- Meta Information

https://aclanthology.org/W17-1603/

# Data Publication

$R^3$ —
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

Community

Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
Call to Action

Data Types Collected

"Medical and Proficiency testing Speech"

"Social Media"

"Anno"

7%

27%

33%

37%

"Written"

R³ —
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

Community

Replicating
and Reporting
Values

Replicating
Results

Responsible
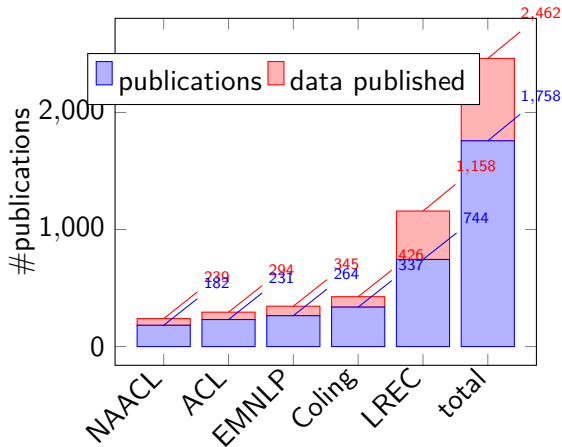
Conclusions &
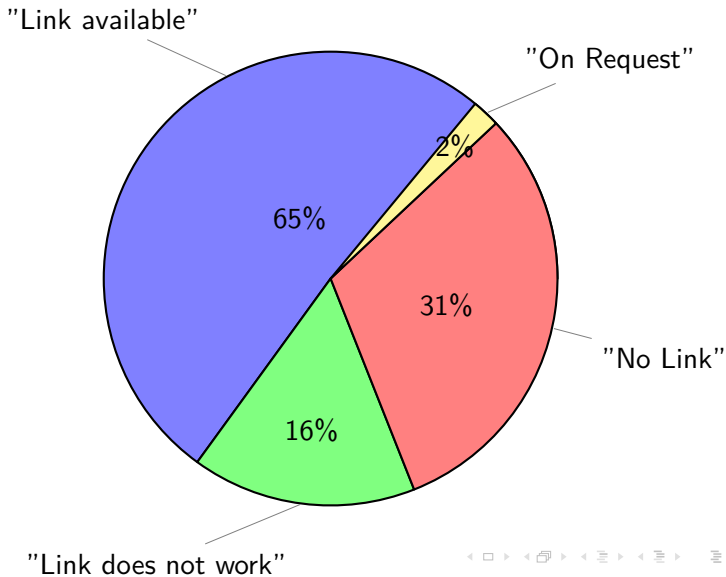Call to Action

# How can I get the data – if at all

$R^3$ –
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

Community

Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
Call to Action

"Link available"

"On Request"

2%

65%

31%

"No Link"

16%

"Link does not work"

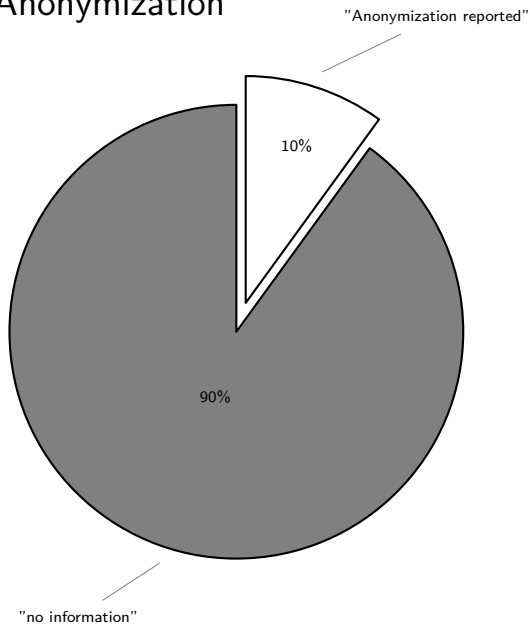# Anonymization

R³ —
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

Community

Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
Call to Action

"Anonymization reported"

10%

90%

"no information"

# Have we evolved?

$R^3$ —
*Responsible Rep*

Margot
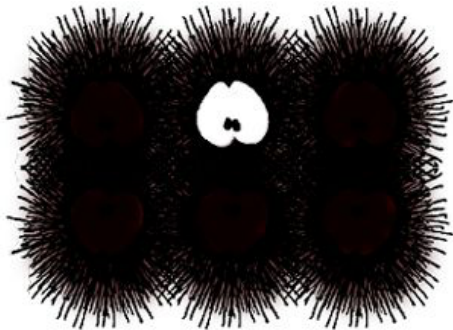Mieskes

Research

**Research
Artifacts**

Community

Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
Call to Action

# Community

$R^3$ —
Responsible Rep

Margot
Mieskes
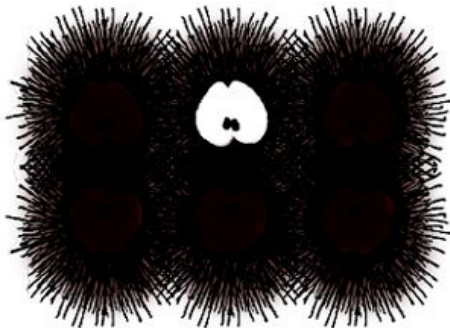
Research

Research
Artifacts

**Community**
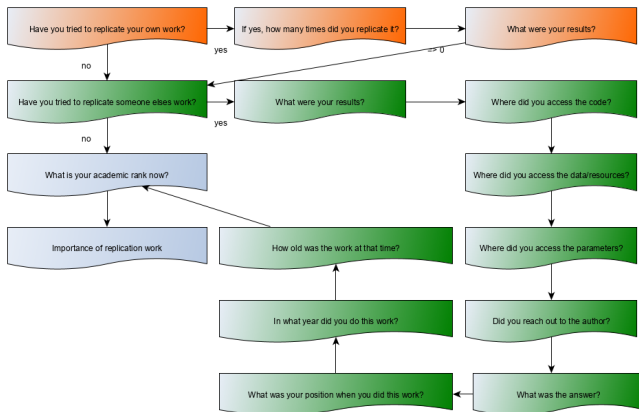
Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
Call to Action

Is it actually a problem?

# What did we do?

R$^3$ –
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

**Community**

Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
Call to Action

# What did we do?

R³ –
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

**Community**

Replicating
and Reporting
Values

Replicating
Results

Responsible

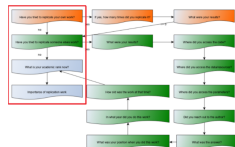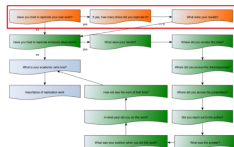Conclusions &
Call to Action

- Have you tried to replicate your own work?
- Have you tried to replicate someone elses work?
- What is your academic rank?
- How important do you rate replication?

https://aclanthology.org/R19-1089/

# What did we do?

- Have you tried to replicate your own work?
- If yes, how many times did you replicate it?
- What were the results?

https://aclanthology.org/R19-1089/

# What did we do?

R³ –
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

**Community**

Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
Call to Action

- Have you tried to replicate someone elses work?
- What were the results?

https://aclanthology.org/R19-1089/

# What did we do?

- Where did you access the code?
- Where did you access the data/resources?
- Where did you access the parameters?

https://aclanthology.org/R19-1089/

# What did we do?

R³ —
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

**Community**

Replicating
and Reporting
Values

Replicating
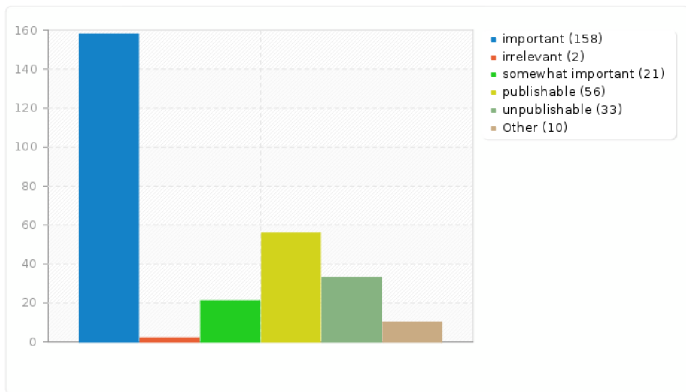Results

Responsible

Conclusions &
Call to Action

- Did you reach out to the author?
- What was the answer?

`https://aclanthology.org/R19-1089/`

# General Stance towards Replication

R³ —
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

Community

Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
Call to Action

Figure: Importance of Replication.

# What are the results?

R³ –
Responsible Rep
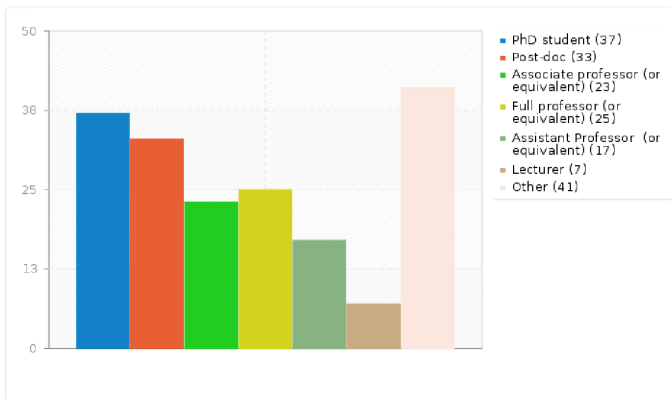
Margot
Mieskes

Research

Research
Artifacts

Community

Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
Call to Action

Part I: Who took part?



Figure: Seniority of Participants.

# What are the results?

R³ —
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

**Community**

Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
Call to Action

Part II: Replicating One's Own Work



Figure: Replicating one's Own Work.

R³ –
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

Community

Replicating
and Reporting
Values

Replicating
Results

Responsible
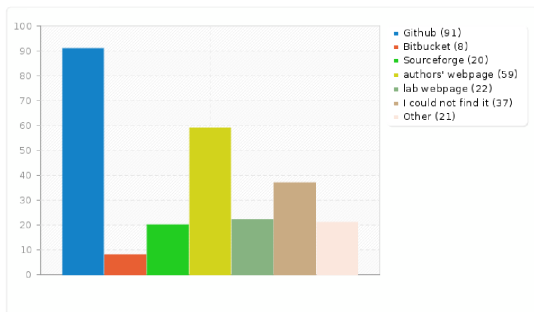
Conclusions &
Call to Action

# What are the results?

And its results:





Figure: Results for Replicating one's Own Work.

# What are the results?

Part II: Replicating Other Researcher's Work



Figure: Replicating Other Researcher's Work.

R³ –
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

**Community**

Replicating
and Reporting
Values

Replicating
Results

Responsible
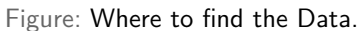
Conclusions &
Call to Action

# What are the results?

And its results:



Figure: Results for Replicating Other Researcher's Work.

# What are the results?

R³ —
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

**Community**

Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
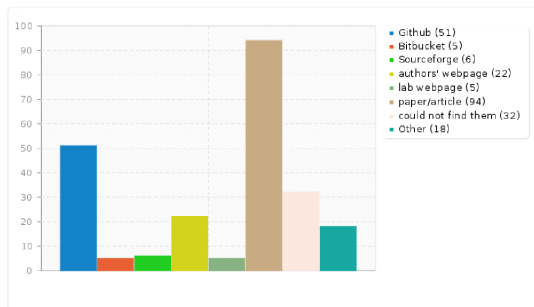Call to Action

Part III: Where are the research artefacts sto
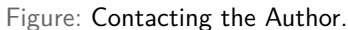


Figure: **Where to find the Code.**

# What are the results?

Part III: Where are the research artefacts sto



Figure: Where to find the Data.

# What are the results?

R³ —
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

Community

Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
Call to Action

Part III: Where are the research artefacts sto



Figure: Where to find the Parameters.

# What are the results?

R³ —
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

Community

Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
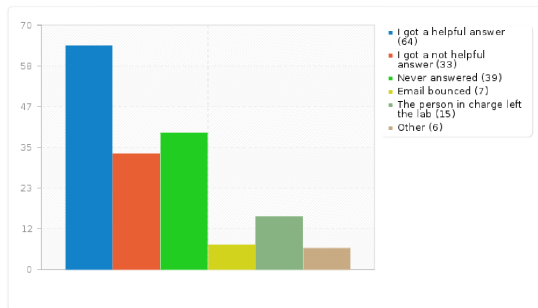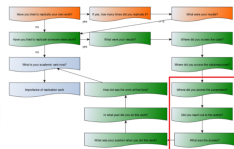Call to Action

Part III: If all fails, contact the author





Figure: **Contacting the Author.**

# What are the results?

R³ —
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

**Community**

Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
Call to Action

Part III: If all fails, contact the author – and



Figure: Types of Answers received from Authors.

R³ –
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

**Community**

Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
Call to Action

## What are our conclusions?

- Replicating one's own work is successful in 23.1% of the answers received
- Replicating others work is successful in 40% of the reported cases – in terms of general conclusions
- Research Artefacts are stored in various places – not all of them with any long-term availability
- Research Artefacts are handed over on a person-to-person basis
- Authors are not generally helpful in answering questions – if they answer at all

:-((((

# Replicating and Reporting Values

What is actually happening?

https://aclanthology.org/2022.insights-1.23/

# Replicating and Reporting Values

What is actually happening?


A use-case from automatic summarization.....


https://aclanthology.org/2022.insights-1.23/

# Background

R³ —
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

Community

Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
Call to Action

# Literature – DUC 2002

$R^3$ –
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

Community

**Replicating
and Reporting
Values**

Replicating
Results

Responsible

Conclusions &
Call to Action

| Loret | Barrera | Mihalcea | official |
|-------|---------|----------|----------|
| S28   | S28     | S27      | S19      |
| S19*  | S21*    | S28      | S28      |
| –     | S29*    | S21      | S21      |
| –     | S23*    | S29*     | S31      |

Table: * did not beat the baseline according to the source paper.

# Results – DUC 2002

R³ –
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

Community

Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
Call to Action

| Citation | S28 | S21 | S19 |
|---|---|---|---|
| Mihalcea(2004) | 0.4703 | 0.4683 | na |
| Mihalcea(2005) | 0.4890 | 0.4869 | na |
| Loret(2010) | 0.4278 | 0.4149 | 0.4082 |
| Barrera(2011) | 0.4781 | 0.4754 | 0.4552 |

# Literature – DUC 2004

R³ –
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

Community

Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
Call to Action

| Citation | ROUGE-1 |
|---|---|
| Original | 0.38224 |
| Yih Wen-Tau (2007)[†] | 0.305 |
| Alguliev (2012)b | 0.3822 |
| Ryang(2012) | 0.3827 |
| Manna(2012)[†] | 0.3913 |
| Rioux(2014)[†] | 0.3828 |
| Ren(2016)[†] | 0.3788 |
| Wang(2017)[†] | 0.3762 |

Table: [†] indicates that parameters have been reported in the publication.

# MEAD – DUC 2004

R³ –
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

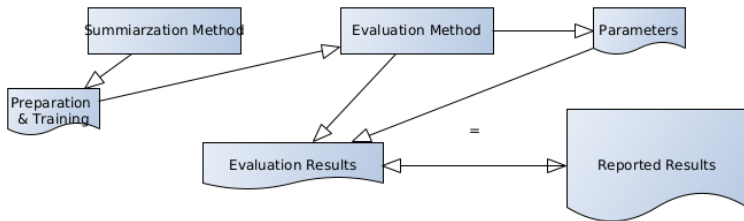Community

Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
Call to Action

| Citation | Result |
|---|---|
| Erkan(2004)a (added features) | 0.38304 |
| Erkan(2004) | 0.3758 |
| Alguliev(2012) | 0.3673 |
| Hong(2014)[†] | 0.3641 |
| re-run | 0.3494 |

Table: [†] indicates that parameters have been reported in the publication.

# SVM

R³ –
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

Community

Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
Call to Action

| Sipos(2012) | re-train & eval (95% CI) |
| --- | --- |
| 0.4066 | 0.3995 (0.3883–0.4117) |

Table: Re-evaluation on DUC 2004 data.

# Take home

$R^3$ —
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

Community

Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
Call to Action

# Replicating Results

Can we at least draw conclusions?

https://aclanthology.org/2023.humeval-1.11/

R³ –
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

Community

Replicating
and Reporting
Values

**Replicating
Results**

Responsible

Conclusions &
Call to Action

# Replicating Results

Can we at least draw conclusions?

An example from the ReproHum Activities

https://aclanthology.org/2023.humeval-1.11/

R³ –
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

Community

Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
Call to Action

## The Original Study

- Presented by Lux and Vu in 2022
- Creating a TTS System with few resources (data & training time)
- Large Multi-Lingual Model
- Fine-Tuning towards target language
- Models
    - Tacotron2
    - FastSpeech2
- Data
    - English, Greek, Spanish, Finnish, Russian, Hungarian, Dutch, French, German
- Target
    - German
    - 30 Minutes of data
    - Training for 2h
- Comparison: 29 hours of German Data exclusively

# Reproducing the Human Evaluation

R³ —
*Responsible Rep*

Margot
Mieskes

Research

Research
Artifacts

Community

Replicating
and Reporting
Values

**Replicating
Results**

Responsible

Conclusions &
Call to Action

# Results

## Comparison

R³ —
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

Community

Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
Call to Action

**FastSpeech2 (Lux and Vu, 2022)**

Proposed
25,3%

No preference
43,4%

Baseline
31,3%

**FastSpeech 2 (Our results)**

Proposed
40,5%

No preference
46,4%

Baseline
13,1%

**Tacotron 2 (Lux and Vu, 2022)**

No preference
37,0%

Proposed
52,0%

Baseline
11,0%

**Tacotron 2 (Our results)**

Proposed
25,7%

No preference
51,8%

Baseline
22,5%

# Reproducing the Technical Background – The Voice

$R^3$ –
Responsible Rep

Margot
Mieskes

Research

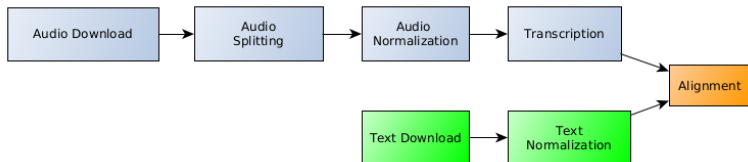Research
Artifacts

Community

Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
Call to Action

Pipeline for creating the Audio-Transcript Data according to Puchtler et al (2021)

# Results

R³ –
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

Community

Replicating
and Reporting
Values

**Replicating
Results**

Responsible

Conclusions &
Call to Action

- Dead Links to Speech Model – fixed by original authors
- Hard-coded links to Textual Data dead – fixable
- Missing packages – fixable
- Faulty files from one speaker – not recoverable

Aborted Replication Attempt

# Reproducing the Technical Background – The Model

$R^3$ –
Responsible Rep

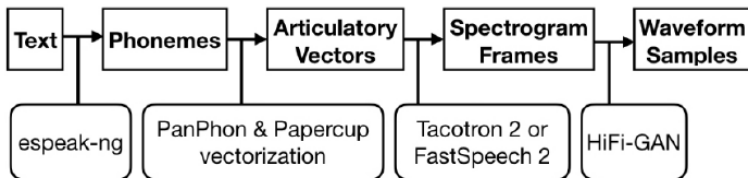Margot
Mieskes

Research

Research
Artifacts

Community

Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
Call to Action

Pipline to create the TTS model according to Lux and Vu (2022)

# Results

| Modell | HW | Dur Prepr | It | T/It | Tot Dur |
|---|---|---|---|---|---|
| Tacotron2 Low Resource | GPU | 1:13 min | 10,020 | 1.25 It/sec | 2:25 hrs |
| Tacotron2 full | GPU | 50:32 min | 100,224 | 1.4 It/sec | 19:54 hrs |
| Tacotron2 Low Resource | CPU | NA | 925 | 22 sec/It | 6 hrs |
| FastSpeech2 Low Resource | GPU | NA | 100,071 | 4.4 It/sec | 6:27 hrs |

- Missing packages – fixed
- FFMPEG issues – fixed
- Two versions of the corpus – fixed?
- Number of Training Iterations – fixed?

# Results

$R^3$ —
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

Community

Replicating
and Reporting
Values

**Replicating
Results**

Responsible

Conclusions &
Call to Action

# But

Notable Differences in

- Sound Quality
- Pronounciation
- Prosody

# Discussion

| Reproduction | Reproducibility | Remarks |
|---|---|---|
| Data set | Reproduction had to be abandoned | Mirrors unavailable, software issues |
| TTS Model | Partially, conclusions were reproduced | Different results, conclusion can be supported |
| Human Evaluation | Values and results not reproducible | Overall conclusion reproducible |

# The Verdict

$R^3$ —
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

Community

Replicating
and Reporting
Values

**Replicating
Results**

Responsible

Conclusions &
Call to Action

- Replication Crisis
- Failure to Replicate Results even in Computer Science
- Failure to find Necessary Research Artefacts in the NLP domain
- Negative Results when Replicating

# Responsible

$R^3$ —
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

Community

Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
Call to Action

So how responsible are you taking your tasks as a researcher?

# Responsible

$R^3$ –
Responsible Rep

Margot
Mieskes

Research

Research
Artifacts

Community
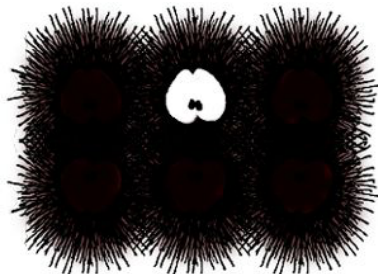
Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
Call to Action

So how responsible are you taking your tasks as a researcher?

- creative
- systematic
- control
- replicate

$R^3 -$
*Responsible Rep*

Margot
Mieskes

Research

Research
Artifacts

Community

Replicating
and Reporting
Values

Replicating
Results

**Responsible**

Conclusions &
Call to Action

# Responsible

So how responsible are you taking your tasks as a researcher?



https://de.wikipedia.org/wiki/Datei:Checkliste.svg

# Checklist or Guidelines?

$R^3$ —
Responsible Rep

Margot
Mieskes

Research

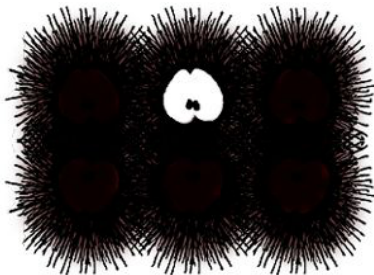Research
Artifacts

Community

Replicating
and Reporting
Values

Replicating
Results

Responsible

Conclusions &
Call to Action

## Questions?

Many thanks to my students Jacob Benz, Jonathan Baum, Christian Stute,
my colleagues Karen Fort, Aurelie Neveol, Kevin B. Cohen,



and thank you for your attention!