# Gender bias in Neural Machine Translation

**Argentina Anna Rescigno**
Eva Vanmassenhove
Johanna Monti
Andy Way

6th October 2020

1

*Proceedings of the 14th Conference of the Association for Machine Translation in the Americas*
*October 6 – 9, 2020, Workshop on the Impact of Machine Translation*
*Page 62*

- **Introduction**
  - o A Note on Terminology
  - o A Quick Problem Sketch

- **Experimental setup**
  - o Compilation of Datasets
  - o Description of the MT systems

- **Results & Analysis**

- **Three main points:**
  - o Why does this kind of bias matter
  - o What is its impact and on whom
  - o Why we need to correct this bias

- **Conclusions and Future Work**

2

*Proceedings of the 14th Conference of the Association for Machine Translation in the Americas*
*October 6 – 9, 2020, Workshop on the Impact of Machine Translation*

*Page 63*

# Introduction

3

*Proceedings of the 14th Conference of the Association for Machine Translation in the Americas*
*October 6 – 9, 2020, Workshop on the Impact of Machine Translation*

*Page 64*

# Introduction: a note on terminology

| **Natural Gender** |
| --- |
| *"Gender based on the **sex** or, for neuter, the lack of sex of the referent of a noun, as English girl (<u>feminine</u>) is referred to by the feminine pronoun she, boy (<u>masculine</u>) by the masculine pronoun he, and table (neuter) by the <u>neuter</u> pronoun it."* <br><br> *Collins Dictionary* 2018, HarperCollins, London, <br> viewed September 2020 <br> http://www.collinsdictionary.com |

4

*Proceedings of the 14th Conference of the Association for Machine Translation in the Americas*
*October 6 – 9, 2020, Workshop on the Impact of Machine Translation*

*Page 65*

| Natural Gender | Grammatical Gender |
|---|---|
| *"Gender based on the **sex** or, for neuter, the lack of sex of the referent of a noun, as English girl (<u>feminine</u>) is referred to by the feminine pronoun she, boy (<u>masculine</u>) by the masculine pronoun he, and table (neuter) by the <u>neuter</u> pronoun it."* | *"Gender based on arbitrary assignment, without regard to the referent of a noun, as in French 'le livre' (masculine), "the book," and German 'das Mädchen' (neuter), "the girl."* |
| *Collins Dictionary* 2018, HarperCollins, London, viewed September 2020 http://www.collinsdictionary.com | *Collins Dictionary* 2018, HarperCollins, London, viewed September 2020 http://www.collinsdictionary.com |

5

*Proceedings of the 14th Conference of the Association for Machine Translation in the Americas*
*October 6 – 9, 2020, Workshop on the Impact of Machine Translation*

*Page 66*

| Natural Gender | Grammatical Gender | Social Gender |
|---|---|---|
| *"Gender based on the **sex** or, for neuter, the lack of sex of the referent of a noun, as English girl (<u>feminine</u>) is referred to by the feminine pronoun she, boy (<u>masculine</u>) by the masculine pronoun he, and table (neuter) by the <u>neuter</u> pronoun it."* | *"Gender based on arbitrary assignment, without regard to the referent of a noun, as in French 'le livre' (masculine), "the book," and German 'das Mädchen' (neuter), "the girl."* | - *Embedded in the lexicon of many languages*<br><br>- *Systematic structural bias.*<br><br>- *Masculine forms the default for generic use.* |
| *Collins Dictionary* 2018, HarperCollins, London, viewed September 2020 http://www.collinsdictionary.com | *Collins Dictionary* 2018, HarperCollins, London, viewed September 2020 http://www.collinsdictionary.com | |

6

**Romance Languages (e.g. ES, FR, IT)**

- animate/persons/animals

↓

grammatical gender = natural gender

- inanimate objects

↓

grammatical gender = arbitrary

*Proceedings of the 14th Conference of the Association for Machine Translation in the Americas*
*October 6 – 9, 2020, Workshop on the Impact of Machine Translation*

*Page 68*

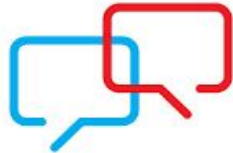| Romance Languages (e.g. ES, FR, IT) | English |
|---|---|
| ● animate/persons/animals<br><br>↓<br><br>grammatical gender = natural gender<br><br><br>● inanimate objects<br><br>↓<br><br>grammatical gender = arbitrary | ● grammatical gender is not inflectional<br><br>● ***pronominal gender*** → gender expressed through the pronouns = natural gender<br><br>● ***gender-neutralization*** of the language |

*Proceedings of the 14th Conference of the Association for Machine Translation in the Americas*
*October 6 – 9, 2020, Workshop on the Impact of Machine Translation*

*Page 69*

**A simple example:**
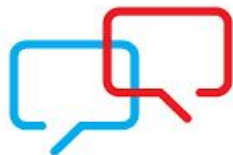
I am happy!

Io sono
content<u>o</u>!

Io sono
content<u>a</u>!

[Natural Gender]

[Grammatical Gender]

I am happy!

Je suis
heureu<u>x</u>!

Je suis
heureu<u>se</u>!

[Natural Gender]

[Grammatical Gender]

9

*Proceedings of the 14th Conference of the Association for Machine Translation in the Americas*
*October 6 – 9, 2020, Workshop on the Impact of Machine Translation*

*Page 70*

| | | Subject gender | Predicative nominative gender | Agreement? |
|---|---|---|---|---|
| **English** | Mark is an efficient <u>nurse</u>. | M | covered | / |
| **Italian** | Mark è <u>un'infermiera</u> efficiente. | M | F | X |
| **French** | Mark est <u>une infirmière</u> efficace. | M | F | X |
| **Spanish** | Mark es <u>una enfermera</u> eficiente. | M | F | X |

*Nov 2019*

➢ **Lack of diversity** → preference for masculine & gender-bias exemptions

➢ **Agreement errors**

10

*Proceedings of the 14th Conference of the Association for Machine Translation in the Americas*
*October 6 – 9, 2020, Workshop on the Impact of Machine Translation*

*Page 71*

# Experimental Setup

11

*Proceedings of the 14th Conference of the Association for Machine Translation in the Americas*
*October 6 – 9, 2020, Workshop on the Impact of Machine Translation*

*Page  72*

## Gender bias in MT

Google Translate

DeepL Translator

Bing Microsoft Translator

- personality adjectives
- profession nouns
- bigender nouns (in Italian)
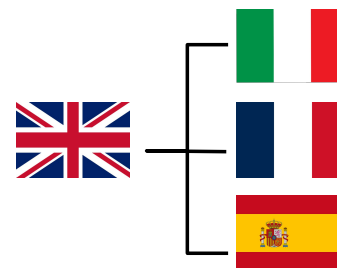  - minimal sentence "I am a(n)..."
  - sentence with a referring adjective

|  | # | Sources |
|---|---|---|
| **Adjectives** | 136 | (I, 2019a); (II, 2019a);(III, 2019) |
| **Professions** | 107 | (I, 2019b); (II, 2019b) |
| **Bigender** | 30 | (Cacciari et al., 1997); (Cacciari et al., 2011) (Thornton and Anna, 2004) |

Table 1: Overview of adjectives, profession and bigender nouns along with the sources from which they were retrieved

2

*Proceedings of the 14th Conference of the Association for Machine Translation in the Americas*
*October 6 – 9, 2020, Workshop on the Impact of Machine Translation*

*Page 73*

# Compilation of Datasets

| | # | Sources |
|---|---|---|
| **Adjectives** | 136 | (I, 2019a); (II, 2019a);(III, 2019) |
| **Professions** | 107 | (I, 2019b); (II, 2019b) |
| **Bigender** | 30 | (Cacciari et al., 1997); (Cacciari et al., 2011) (Thornton and Anna, 2004) |

Table 1: Overview of adjectives, profession and bigender nouns along with the sources from which they were retrieved

| English | Italian | | French | | Spanish | |
|---|---|---|---|---|---|---|
| I am an assistant. | Sono un assistente. | M | Je suis un assistant. | M | Soy asistente. | * |
| I am a beautiful assistant. | **Sono una bellissima assistente.** | **F** | **Je suis une belle assistante.** | **F** | **Soy una bella asistente.** | **F** |
| I am an efficient assistant. | Sono un assistente efficiente. | M | Je suis un assistant efficace. | M | Soy un asistente eficiente. | M |

| | | | | | | |
|---|---|---|---|---|---|---|
| I am a translator. | Sono un traduttore. | M | Je suis un traducteur. | M | Soy un traductor. | M |
| I am a beautiful translator. | **Sono una bellissima traduttrice.** | **F** | **Je suis une belle traductrice.** | **F** | **Soy una bella traductora.** | **F** |
| I am an efficient translator. | Sono un traduttore efficiente. | M | Je suis un traducteur efficace. | M | Soy un traductor eficiente. | M |

13

*Proceedings of the 14th Conference of the Association for Machine Translation in the Americas*
*October 6 – 9, 2020, Workshop on the Impact of Machine Translation*

*Page 74*

**Google Translate**

- 2003

- statistical MT system

- 2016 →  neural MT system

- 2018 →  double alternatives on word level

14

*Proceedings of the 14th Conference of the Association for Machine Translation in the Americas*
*October 6 – 9, 2020, Workshop on the Impact of Machine Translation*

*Page 75*

Google Translate

**DeepL Translator**

- 2017
- convolutional neural networks
- Linguee database (dictionary)
- nine languages supported
- provides not morphological alternatives
- serves also as glossary

15

*Proceedings of the 14th Conference of the Association for Machine Translation in the Americas*
*October 6 – 9, 2020, Workshop on the Impact of Machine Translation*

*Page 76*

Google Translate

DeepL Translator

**Bing Microsoft Translator**

- originally a statistical MT system

- switched to a neural system

- does not provides alternatives but

- provides examples of usage

16

*Proceedings of the 14th Conference of the Association for Machine Translation in the Americas*
*October 6 – 9, 2020, Workshop on the Impact of Machine Translation*

*Page 77*

# Results & Analysis

17

*Proceedings of the 14th Conference of the Association for Machine Translation in the Americas*
*October 6 – 9, 2020, Workshop on the Impact of Machine Translation*

*Page 78*

❏   ADJECTIVES

| ADJ | GT | BMT | DL |
|---|---|---|---|
| F | 37.3 | 1.5 | 22.8 |
| M | **39.2** | **58.8** | **45.6** |
| N | 20.7 | 33.1 | 26.5 |
| Other | 2.8 | 6.5 | 5.1 |
| Total | 100 | 100 | 100 |

Table 2: Results in % for male (M), female (F) and neutral (N) adjectives generated for EN → IT for GT, BMT and DL. The "Other" label includes all results obtained that do not correspond to the "adjective" category

18

*Proceedings of the 14th Conference of the Association for Machine Translation in the Americas*
*October 6 – 9, 2020, Workshop on the Impact of Machine Translation*

*Page 79*

❏ NOUNS

| NOUN | GT | BMT | DL |
|------|------|------|------|
| F | 35.8 | 0.9 | 7.5 |
| M | **46.1** | **60.4** | **60.4** |
| N | 17.6 | 28.3 | 28.3 |
| Other | 0.6 | 10.5 | 3.7 |
| Total | 100 | 100 | 100 |

Table 3: Results in % for male (M), female (F) and neutral (N) nouns generated for EN → IT for GT, BMT and DL. The "Other" label includes all results obtained that do not correspond to the "noun" category

19

*Proceedings of the 14th Conference of the Association for Machine Translation in the Americas*
*October 6 – 9, 2020, Workshop on the Impact of Machine Translation*

*Page 80*

| BMT | IT | | | FR | | | ES | | |
|---|---|---|---|---|---|---|---|---|---|
| | F | M | N | F | M | N | F | M | N |
| no adj. | 10.0 | **86.7** | $Q^*$ | 10.0 | **63.3** | 26.7 | 3.3 | **66.7** | 30.0 |
| beautiful | **63.3** | 36.7 | 0.0 | 43.3 | **56.7** | 0.0 | **66.7** | 33.3 | 0.0 |
| other adj. | 13.3 | **83.3** | $Q^*$ | 3.3 | **96.7** | 0.0 | 6.7 | **93.3** | 0.0 |
| **DL** | **IT** | | | **FR** | | | **ES** | | |
| | F | M | N | F | M | N | F | M | N |
| no adj | 30.0 | **70.0** | 0.0 | 20.0 | **63.3** | 16.7 | 3.3 | **76.6** | 20.0 |
| beautiful | **83.3** | 16.7 | 0.0 | **73.3** | 26.7 | 0.0 | **96.7** | 3.3 | 0.0 |
| other adj. | **53.3** | 43.3 | $Q^*$ | 13.3 | **83.3** | 3.3 | 6.7 | **93.3** | 0.0 |
| **GT** | **IT** | | | **FR** | | | **ES** | | |
| | F | M | N | F | M | N | F | M | N |
| no adj. | 6.7 | **93.3** | 0.0 | 6.7 | **90.0** | 3.3 | 3.3 | **66.7** | 30.0 |
| beautiful | 43.3 | **56.7** | 0.0 | **80.** | 20.0 | 0.0 | **80.0** | 20.0 | 0.0 |
| other adj. | 3.3 | **96.7** | 0.0 | 3.3 | **96.7** | 0.0 | 3.3 | **96.7** | 0.0 |

Table 4: Results in % for male (M), female (F) and neutral (N) forms generated for EN → IT, FR and ES for BMT, DL and GT

- *beautiful*

other adjectives:

- *efficient*
- *intelligent*
- *sad*
- *famous*

20

*Proceedings of the 14th Conference of the Association for Machine Translation in the Americas*
*October 6 – 9, 2020, Workshop on the Impact of Machine Translation*

*Page 81*

| BMT | IT | | | FR | | | ES | | |
|---|---|---|---|---|---|---|---|---|---|
| | F | M | N | F | M | N | F | M | N |
| no adj. | 10.0 | **86.7** | $Q^*$ | 10.0 | **63.3** | 26.7 | 3.3 | **66.7** | 30.0 |
| beautiful | **63.3** | 36.7 | 0.0 | 43.3 | **56.7** | 0.0 | **66.7** | 33.3 | 0.0 |
| other adj. | 13.3 | **83.3** | $Q^*$ | 3.3 | **96.7** | 0.0 | 6.7 | **93.3** | 0.0 |
| DL | IT | | | FR | | | ES | | |
| | F | M | N | F | M | N | F | M | N |
| no adj | 30.0 | **70.0** | 0.0 | 20.0 | **63.3** | 16.7 | 3.3 | **76.6** | 20.0 |
| beautiful | **83.3** | 16.7 | 0.0 | **73.3** | 26.7 | 0.0 | **96.7** | 3.3 | 0.0 |
| other adj. | **53.3** | 43.3 | $Q^*$ | 13.3 | **83.3** | 3.3 | 6.7 | **93.3** | 0.0 |
| GT | IT | | | FR | | | ES | | |
| | F | M | N | F | M | N | F | M | N |
| no adj. | 6.7 | **93.3** | 0.0 | 6.7 | **90.0** | 3.3 | 3.3 | **66.7** | 30.0 |
| beautiful | 43.3 | **56.7** | 0.0 | **80.** | 20.0 | 0.0 | **80.0** | 20.0 | 0.0 |
| other adj. | 3.3 | **96.7** | 0.0 | 3.3 | **96.7** | 0.0 | 3.3 | **96.7** | 0.0 |

Table 4: Results in % for male (M), female (F) and neutral (N) forms generated for EN → IT, FR and ES for BMT, DL and GT
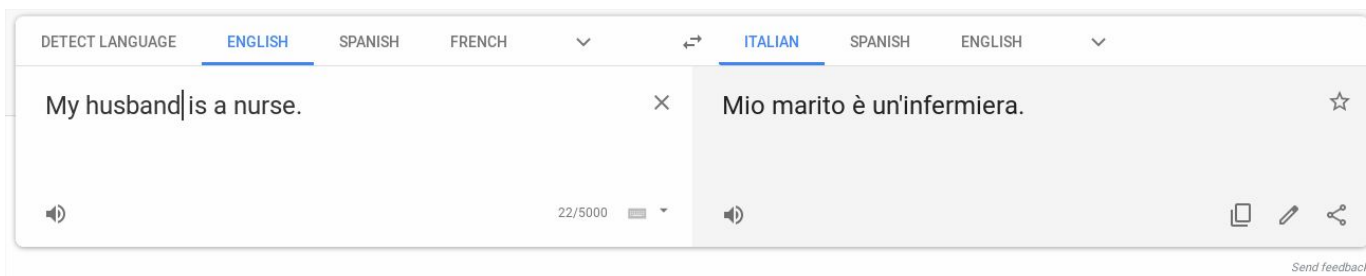
- *beautiful*

other adjectives:

- *efficient*
- *intelligent*
- *sad*
- *famous*

21

*Proceedings of the 14th Conference of the Association for Machine Translation in the Americas*
*October 6 – 9, 2020, Workshop on the Impact of Machine Translation*

*Page 82*

# iMpacT

22

*Proceedings of the 14th Conference of the Association for Machine Translation in the Americas*
*October 6 – 9, 2020, Workshop on the Impact of Machine Translation*
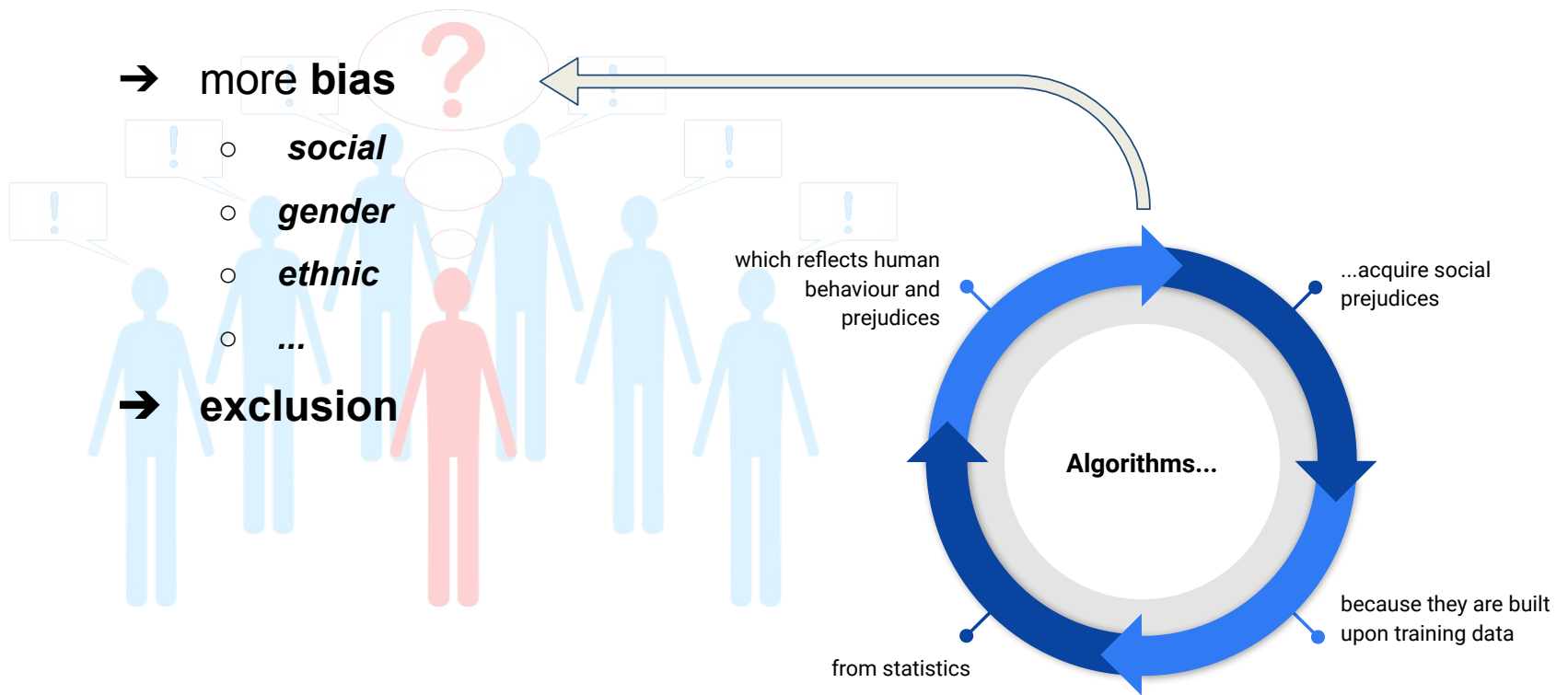
*Page 83*

- **From a linguistic point of view:**
  - Avoiding basic gender agreement mistakes



- **From a technological point of view:**
  - Solving these issues is not trivial (see attempts Google)

  - Black box of NLP (we have no/little control over the actual output that are being generated)

- **From a societal/ethical point of view:**
  - Identifying biases in current state-of-the-art systems is important so they don't end up getting mistaken for 'objective' translations

  - if an MT system is being used without human in the loop: real-world consequences

23

*Proceedings of the 14th Conference of the Association for Machine Translation in the Americas*
*October 6 – 9, 2020, Workshop on the Impact of Machine Translation*

*Page 84*

**Break the cycle**

➔ more **bias**
  ○ *social*
  ○ *gender*
  ○ *ethnic*
  ○ *...*

➔ **exclusion**

which reflects human behaviour and prejudices

...acquire social prejudices

**Algorithms...**

because they are built upon training data

from statistics

24

*Proceedings of the 14th Conference of the Association for Machine Translation in the Americas*
*October 6 – 9, 2020, Workshop on the Impact of Machine Translation*

*Page 85*

# Conclusion and
## Future Work

25

*Proceedings of the 14th Conference of the Association for Machine Translation in the Americas*
*October 6 – 9, 2020, Workshop on the Impact of Machine Translation*

*Page 86*

## Conclusion:

- Remove gender bias in training data

- Train algorithms to address the problem

- Stop using masculine "neutral" in machine learning texts

- Evaluation of gender phenomena is challenging

## Future Work:

- Extend to other language pairs (different languages → different gender phenomena)

- Larger evaluation of more diverse set of words

- Create language specific challenge sets to evaluate how biased is an MT system

- Train our own MT system to verify whether machine bias influences the output of the translation

26

*Proceedings of the 14th Conference of the Association for Machine Translation in the Americas*
*October 6 – 9, 2020, Workshop on the Impact of Machine Translation*

*Page 87*

# Thank you for your attention!

27

Bond, E., 2020. Cambridge Researchers Tackle Neural Machine Translation'S Gender Bias | Slator. [online] Slator. Available at: <https://slator.com/machine-translation/cambridge-researchers-tackle-neural-machine-translations-gender-bias/> [Accessed 29 September 2020].

Caliskan, A., Bryson, J. and Narayanan, A., 2017. Semantics derived automatically from language corpora contain human-like biases. Science, 356(6334), pp.183-186.

Devlin, H., 2017. AI Programs Exhibit Racial And Gender Biases, Research Reveals. [online] the Guardian. Available at: <https://www.theguardian.com/technology/2017/apr/13/ai-programs-exhibit-racist-and-sexist-biases-research-reveals> [Accessed 29 September 2020].

Monti, J., 2017. Questioni di genere in traduzione automatica. In: A. de Meo, L. di Pace, A. Manco and J. Monti, ed., AI femminile. Scritti linguistici in onore di Cristina Vallini. Firenze: Cesati, pp.411-431.

Prates, M., Avelar, P. and Lamb, L., 2019. Assessing gender bias in machine translation: a case study with Google Translate. Neural Computing and Applications, 32(10), pp.1-19.

Saunders, D. and Byrne, B., 2020. Reducing Gender Bias in Neural Machine Translation as a Domain Adaptation Problem. *arXiv: 2004.04498v3*

Vanmassenhove, E., Shterionov, D. and Way, A., 2019. Lost in Translation: Loss and Decay of Linguistic Richness in Machine Translation. *arXiv: 1906.12068*

Zou, J. and Schiebinger, L., 2018. AI can be sexist and racist — it's time to make it fair. Nature, 559(7714), pp.324-326.

28

*Proceedings of the 14th Conference of the Association for Machine Translation in the Americas*
*October 6 – 9, 2020, Workshop on the Impact of Machine Translation*

*Page 89*

# Contact info

Argentina A. Rescigno: argentina.res@gmail.com

Eva Vanmassenhove: vanmassenhove.eva@gmail.com

Johanna Monti: johmonti@gmail.com

Andy Way: andy.way@adaptcentre.ie

29

*Proceedings of the 14th Conference of the Association for Machine Translation in the Americas*
*October 6 – 9, 2020, Workshop on the Impact of Machine Translation*

*Page 90*