# Mitigating Negative Transfer with Task Awareness for Sexism, Hate Speech, and Toxic Language Detection

Hate Speech Detection

©2020 Garth German garthtoons.com 🚯 🕑 💿 @garthtoons

[YES] "Shut up you motherfucker. What did

**[NO]** "This is inhumane Karma is a bitch she

get around these brainless heartless assholes!

you do for politics... nothing"

REALIN' Surroundings "German

Example:

OFF TO CHECK YOUR

Angel Felipe Magnossão de Paula<sup>69</sup><sup>69</sup>, Paolo Rosso<sup>69</sup> and Damiano Spina<sup>69</sup>

Department of Computer Systems and Computation, Universitat Politècnica de València

School of Computing Technologies, RMIT University

# **BIG PICTURE**

In the field of machine learning, the common strategy is to apply the Single-Task Learning approach in order to train a supervised model to solve a specific task. Training a robust model requires a lot of data and a significant amount of computational resources, making this solution unfeasible in cases where data are unavailable or expensive to gather. Therefore another solution, based on the sharing of information between tasks, has been developed: Multi-task Learning (MTL). Despite the recent developments regarding MTL, the problem of negative transfer has still to be solved. Hence, we propose a new approach to mitigate the negative transfer problem based on the task awareness concept. The proposed approach results in diminishing the negative transfer together with an improvement of performance over classic MTL solution.





Text

### **Example:**



UNIVERSITAT

POLITÈCNICA

DE VALÈNCIA

UNIVERSITY

 $0.980 \pm 0.003$ 

 $0.988 \pm 0.003$ 

HatEval-2019

F1-macro

 $0.730 \pm 0.022$ 

 $0.764 \pm 0.021$ 

 $0.778 \pm 0.020$ 

 $0.777 \pm 0.020$ 

 $0.773 \pm 0.021$ 

 $0.789 \pm 0.020$ 

**0.790** ± 0.020

 $0.782 \pm 0.020$ 

 $0.786 \pm 0.020$ 

 $0.786 \pm 0.020$ 

 $0.789 \pm 0.020$ 

# PROBLEM



**Toxic Language Detection** 

Example: [YES] "Woman behind the wheel, be careful!" [NO] "not all white women are like that"

Example: **[YES]** "To these fucking mangy animals lock them up and throw away the key" [NO] "Against misinformation data"



Snippet Mechanism 2: Task Embedding (TE)

# **EXPERIMENTS**

- **Cross-validation**
- Train-test split



#### TABLE IV **RESULTS OF THE CROSS-VALIDATION EXPERIMENT WITH 95% CONFIDENCE INTERVALS**

Task Identification

Vector

Model	Task Heads	EXIST-2021	DETOXIS-2021	HatEval-2019
		Accuracy	F1-score	F1-macro
STL	Sexism	$0.789 \pm 0.011$		
	Toxic-language	. —	$0.640 \pm 0.014$	-
	Hate-speech	—	-	$0.846 \pm 0.009$
MTL	Sexism + Toxic-language	$0.788 \pm 0.011$	$0.628 \pm 0.014$	_
	Sexism + Hate-speech	$0.791 \pm 0.011$	_	$0.843 \pm 0.009$
	Toxic-language + Hate-speech	-	$0.632 \pm 0.014$	$0.841 \pm 0.009$
	Toxic-language + Hate-speech + Sexism	$0.799 \pm 0.010$	$0.634 \pm 0.014$	$0.842 \pm 0.009$
MTL-TAI	Sexism + Toxic-language	$0.799 \pm 0.010$	$0.649 \pm 0.014$	-
	Sexism + Hate-speech	$0.805 \pm 0.010$	_	$0.984 \pm 0.003$
	Toxic-language + Hate-speech	_	$0.649 \pm 0.014$	$0.988 \pm 0.003$
	Toxic-language + Hate-speech + Sexism	$0.800 \pm 0.010$	$0.650 \pm 0.014$	$0.980 \pm 0.003$
	Sexism + Toxic-language	$0.797 \pm 0.011$	$0.653 \pm 0.014$	
	Sexism + Hate-speech	$0.806 \pm 0.010$	_	<b>0.992</b> ± 0.002

CONTRIBUTION

## **OUR SOLUTION**

in a drop in performance.

Task Awareness (TA)

- Model aware of task
- LT based on the task

#### Example:

<CL> Tweet <SEP> Sexism Detection <

The work of Paolo Rosso was in the framework of the FairTransNLP-Stereotypes research project (PID2021-124361OB-C31) on Fairness and Transparency for equitable NLP applications in social media: Identifying stereotypes and prejudices and developing equitable systems, funded by MCIN/AEI/10.13039/501100011033 and by ERDF, EU A way of making Europe. Damiano Spina is the recipient of an Australian Research Council DECRA Research Fellowship (DE200100064).

Model 1: MTL Task-Aware Input (MTL-TAI)

Output 1



Mechanism 1: Task-Aware Input (TAI)

(1) We propose the use of the Task Awareness (TA) concept to mitigate the negative transfer problem during MTL training. (2) Design of the Task-Aware Input (TAI) mechanism to grant the MTL models with task awareness ability to mitigate negative transfer and even improve results compared with traditional MTL models. (3) Design of the Task Embedding (TE) mechanism to give MTL models task recognition capability to diminish negative transfer and improve the results over classic MTL solutions. (4) Creation and validation of two unified architectures to detect Sexism, Hate Speech, and Toxic Language in text comments. (5) Our proposed method outperforms the State-Of-The-Art on two public benchmarks for Sexism and Hate Speech detection: (i) EXIST-2021 and (ii) HatEval-2019 datasets.

