

# DIKé project: some attempts to improve fairness in (compressed) language models

Julien Velcin

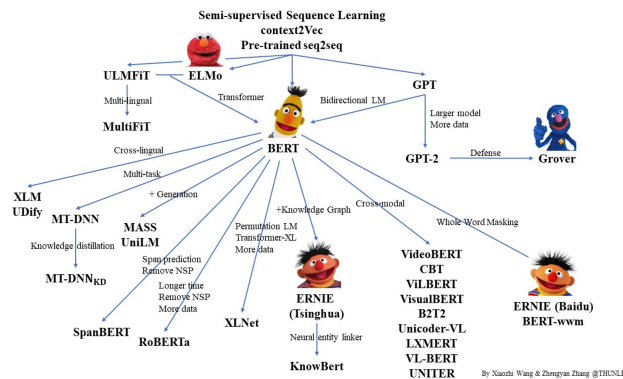
Laboratoire ERIC – Université Lumière Lyon 2

<https://eric.univ-lyon2.fr/jvelcin/>

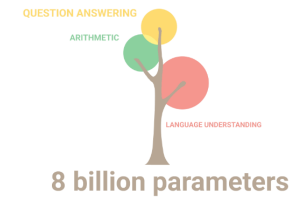
## Outline

- Introduction to LLMs, compression and fairness
- Presentation of the DIKé project
- Focus on hate speech detection (work of I. Proskurina)
- Conclusion and ongoing research on LLMs

## From pre-trained Language Models...



## ...to Large Language Models (LLMs)



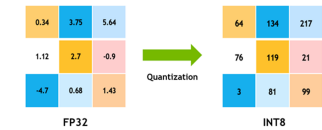
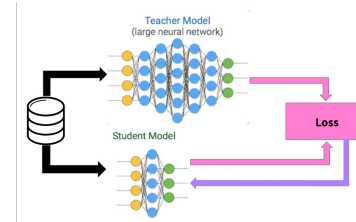
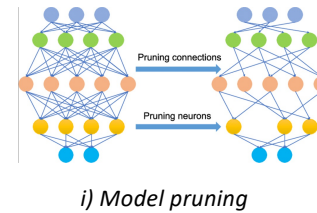
## Model compression (1)

- Transformers are **over-parametrized** (#parameters >> #data)
- Can we just **train smaller** Transformers?  
No (lottery ticket hypothesis)
- **Mobile** devices

Model	# Parameters
BERT	110M
BERT-large	340M
GPT-2	1.5B
LLaMA	65B
GALACTICA	120B
GPT-3	175B
PaLM	540B

5

## Model compression (2)



6

## Bias and Fairness

- **Bias**: discrepancy between the correct way of reasoning, which ensures the validity of the conclusions we draw, and the actual process of reasoning. Often, biases arise from the application of heuristics. Some biases can be related to identity features, such as gender or religion.
- **Fairness** in AI: how to correct algorithmic bias in automated decision processes based on machine learning models, *in particular when the biases target groups of people* (e.g., christians).

7

## DIKé project <https://www.anr-dike.fr>

- Funded by ANR (AAPG 2021, 2022-2025)
- Partners
  - Laboratoire Hubert Curien (LabHC), Université Jean Monnet
  - Laboratoire ERIC, Université Lumière Lyon 2
  - Naver Labs
- Expectations
  - evaluation framework and methodology for evaluating fairness of NLP systems
  - English but also French datasets for fairness and ethics of NLP systems
  - new compressed, fairer language models

8

## Partenaires (PRCE)



Christophe GRAVIER (PR)  
François JACQUENET (PR)  
Antoine GOURRU (MCF)  
Thibaud LETENO (PHD)

Charlotte LACLAU (MCF)  
(Télécom Paris)



Julien VELCIN (PR)  
Guillaume METLZER (MCF)  
Adrien GUILLE (MCF)  
Irina Proskurina (PHD)



Vassilina NIKOULINA (R. Sc)  
Caroline BRUN (Sr Sc.)

9

## Some recent contribution on LLMs

- Naver Labs
  - SMA<sup>LL</sup>-100: Introducing Shallow Multilingual Machine Translation Model for Low-Resource Languages (EMNLP 2023)
  - What Do Compressed Multilingual Machine Translation Models Forget? (Findings of EMNLP 2023)
- LabHC
  - An Investigation of Structures Responsible for Gender Bias in BERT and DistilBERT (IDA 2023)
  - Fair Text Classification with Wasserstein Independence (EMNLP 2024)
- ERIC
  - **The Other Side of Compression: Measuring Bias in Pruned Transformers** (IDA 2023)
  - Mini Minds: Exploring Bebeshka and Zlata Baby Models (CoNLL 2023)
  - When Quantization Affects Confidence of Large Language Models? (NAACL Findings 2024)

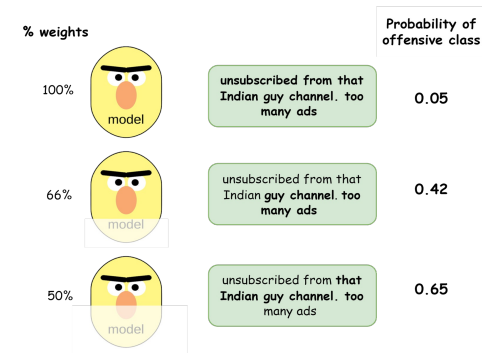
10

## The Other Side of Compression: Measuring Bias in Pruned Transformers (IDA 2023)

- Work of **Irina Proskurina** (PhD student with G. Metzler and me)
- We measure identity-based **bias** in pruned Transformer LMs
- We study **which group** of encoder **layers** (bottom, middle or upper) can be efficiently pruned without biased outcomes
- We propose **word-level supervision** in pruned Transformer LMs as a debiasing method

11

## Bias in Hate Speech Classification

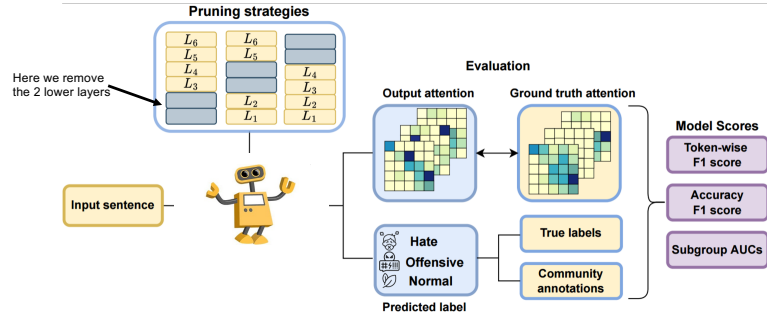


Bias = Compressed LM classifies neutral text as offensive and pays 'attention' to sensitive attributes

13

# Methodology

- 1) Prune Transformer (e.g., BERT)
- 2) Fine-tune Transformer on hate speech classification task
- 3) Evaluate performance, bias, and explainability of fine-tuned pruned Transformers

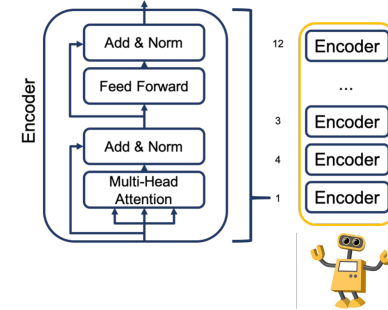


14

## Step 1) Transformers With Pruned Layers (1)

Pruning  $K$  Layers From Transformer:

- Upper = {12,11,10...}
- Bottom = {1,2,3,...}
- Symmetric = {6,7}
- Alternate Odd = {11,9, 7, ...}
- Alternate Even = {12,10,8...}
- Contribution\*-based



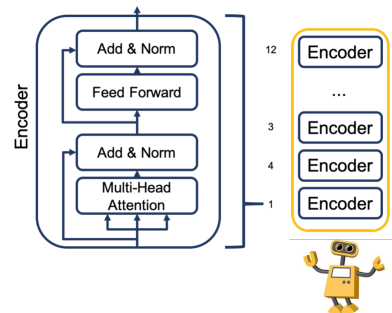
\*Measured with cosine similarity of hidden states:  $\varphi_s(l) = \cos(Z_{l-1}, Z_l)$

15

## Step 1) Transformers With Pruned Layers (2)

Pruning  $K$  Layers From Transformer:

- LMs: BERT, RoBERTa, DistilBERT, DistilRoBERTa
- $K = \{2,4,6\}$  for BERT-based LMs
- $K = \{1,2,3\}$  for Distilled LMs
- Contribution\*-based strategy:
  - BERT: {5, 10, 9, 7, 2, 4}
  - RoBERTa: {1, 2, 6, 8, 9, 4}
  - DistilBERT: {2, 3, 4}
  - DistilRoBERTa: {6, 2, 3}

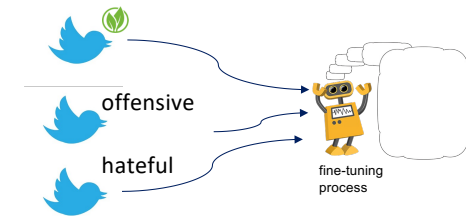


\*Measured with cosine similarity of hidden states:  $\varphi_s(l) = \cos(Z_{l-1}, Z_l)$

16

## Step 2) Fine-tune Transformer on hate speech classification task

- Given input text, Transformer is fine-tuned to classify text as offensive, hateful or neutral
- Data: HATEXPLAIN
- Fine-tuning: Cross-Entropy loss  $L(\theta) = -\sum_{i=1}^c y_i \log(\hat{y}_i)$ , where  $i \in \{\text{offensive, hateful, neutral}\}$ ,  $y_i$  is true class,  $\hat{y}_i$  output model probability of  $i$  belongs to one of the classes



17

### Step 3) Evaluate Bias in Compressed models (1)

- Data divided in 4 domains:  $D_t^+, D_t^-, D_{\setminus t}^+, D_{\setminus t}^-$
- $t$  – (hate) target community, +/- : class, neutral or hateful
- Subgroup AUC =  $AUC(D_t^+, D_t^-)$
- Background Positive Subgroup Negative: BPSN =  $AUC(D_{\setminus t}^+, D_{\setminus t}^-)$
- Background Negative Subgroup Positive: BNSP =  $AUC(D_{\setminus t}^-, D_{\setminus t}^+)$

Subgroup AUC

d	pred	GT
12	0.76	hateful
2	0.76	neutral
307	0.73	neutral
17	0.72	hateful
1	0.69	neutral
101	0.61	neutral
...		

documents related to the group t

18

### Step 3) Evaluate Explainability in Compressed models (2)

- Explainability performance: token-level F1-score
- F1-score(output attention, ground truth attention)
- Output attentions=top-5 tokens with higher attention weights
- Ground truth attention = annotations in red

<user>: I got a guilty pleasure and it is country music and hillbilly movies and tv shows about rednecks hunting in the woods... trailer trash

poor plump thing

<sup>a</sup>Annotator 1: Target labels: *Economic, Caucasian*  
<sup>b</sup>Annotator 2: Target labels: *Economic*  
<sup>c</sup>Annotator 3: Target labels: *Caucasian*

**Rationales**  
= words targeting a given group

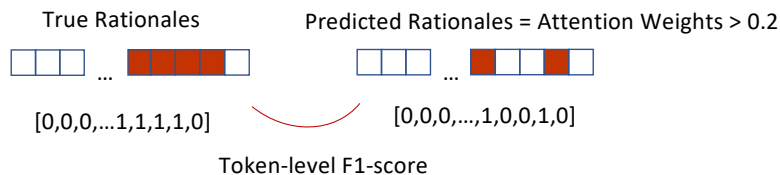
19

### Step 3) Evaluate Explainability in Compressed models (3)

<user>: I got a guilty pleasure and it is country music and hillbilly movies and tv shows about rednecks hunting in the woods... trailer trash

poor plump thing

<sup>a</sup>Annotator 1: Target labels: *Economic, Caucasian*  
<sup>b</sup>Annotator 2: Target labels: *Economic*  
<sup>c</sup>Annotator 3: Target labels: *Caucasian*



20

### Step 3) Evaluate Bias in Compressed models (4)

If the impact of compression is uniform, then the shift in scores achieved on the texts mentioning a target community  $t$  should also be uniform compared to the overall scores shift. That forms our null hypothesis  $H_0$ .

$$H_0 : \beta_0^t - \beta_0 = \beta_c^t - \beta_c \leftarrow \text{no significant difference}$$

$$H_1 : \beta_0^t - \beta_0 \neq \beta_c^t - \beta_c \leftarrow \text{significant difference}$$

21

## Results: Compressed LMs are prone to bias

Model	full model		4 layers removed				
	Layers	F1 score	Token F1 score	Count Subgroup	Signif BNSP	Target BPSN	Classes
BERT	12/12	67.28±0.13	48.58±.28	-	-	-	-
	10/12	65.31±0.17	38.35±4.11	2	0	1	1
	8/12	64.82±0.15	32.57±4.06	2	0	2	2
	6/12	63.46±0.21	34.4±3.87	4	0	2	2
DistilBERT	6/6	66.19±0.44	43.31±3.42	-	-	-	-
	5/6	66.08±0.62	42.77±4.13	0	0	0	0
	4/6	65.66±0.51	42.1±3.98	3	0	1	1
	3/6	64.31±0.83	39.81±4.22	3	1	2	2
RoBERTa	12/12	83.42±0.4	46.64±3.51	-	-	-	-
	10/12	81.46±0.41	39.37±4.61	4	2	2	2
	8/12	78.67±0.58	38.49±4.23	6	3	4	4
	6/12	77.08±0.33	24.47±4.08	6	5	5	5
DistilRoBERTa	6/6	82.02±0.36	42.08±5.24	-	-	-	-
	5/6	81.08±0.4	33.2±4.75	3	0	2	2
	4/6	77.06±0.48	32.76±5.21	3	2	4	4
	3/6	74.05±0.43	32.6±4.61	6	5	6	6

Performance of original and pruned models on HATEXPLAIN test set

22

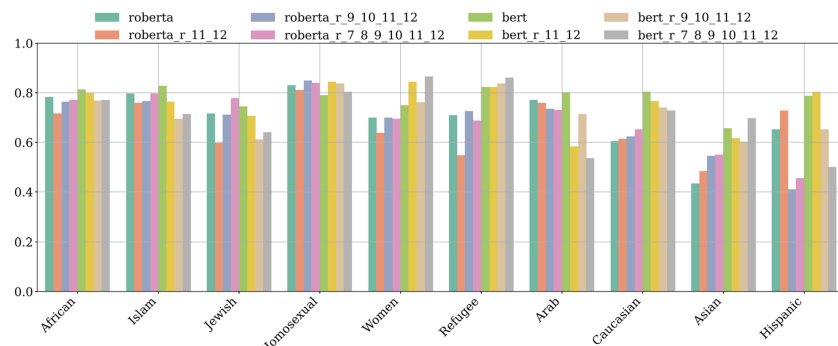
## Results: Compressed LMs rely on unimportant tokens

Model	Layers	F1 score	Token F1 score	Count Subgroup	Signif BNSP	Target BPSN	Classes
BERT	12/12	67.28±0.13	48.58±.28	-	-	-	-
	10/12	65.31±0.17	38.35±4.11	2	0	1	1
	8/12	64.82±0.15	32.57±4.06	2	0	2	2
	6/12	63.46±0.21	34.4±3.87	4	0	2	2
DistilBERT	6/6	66.19±0.44	43.31±3.42	-	-	-	-
	5/6	66.08±0.62	42.77±4.13	0	0	0	0
	4/6	65.66±0.51	42.1±3.98	3	0	1	1
	3/6	64.31±0.83	39.81±4.22	3	1	2	2
RoBERTa	12/12	83.42±0.4	46.64±3.51	-	-	-	-
	10/12	81.46±0.41	39.37±4.61	4	2	2	2
	8/12	78.67±0.58	38.49±4.23	6	3	4	4
	6/12	77.08±0.33	24.47±4.08	6	5	5	5
DistilRoBERTa	6/6	82.02±0.36	42.08±5.24	-	-	-	-
	5/6	81.08±0.4	33.2±4.75	3	0	2	2
	4/6	77.06±0.48	32.76±5.21	3	2	4	4
	3/6	74.05±0.43	32.6±4.61	6	5	6	6

Performance of original and pruned models on HATEXPLAIN test set

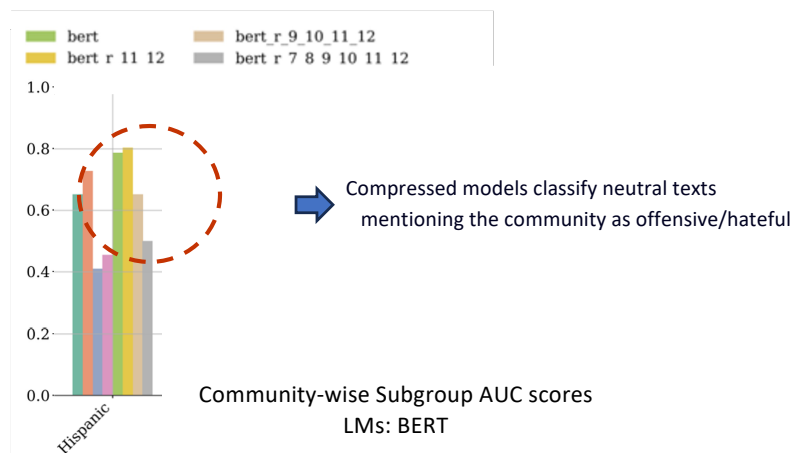
23

## Results: The impact of compression is not uniform



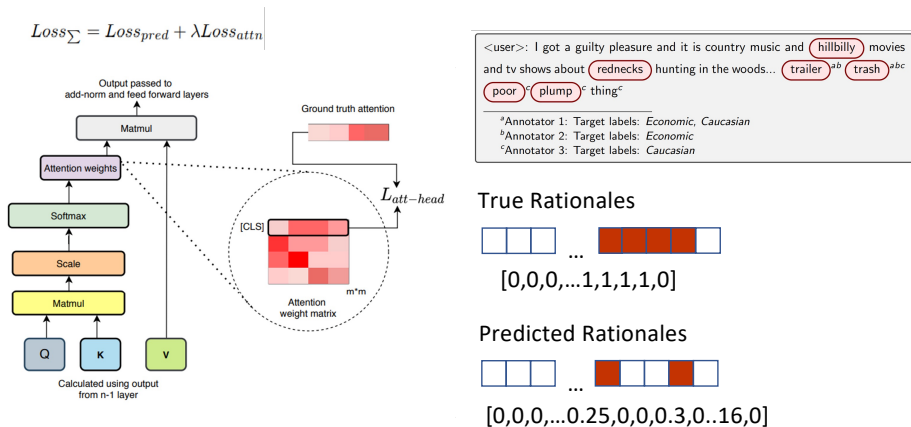
24

## Results: The impact of compression is not uniform



25

## Solution: Supervised Attention learning



26

## Results: Fine-tuning with attention loss compensates for fairness loss

Model	$\lambda$	F1 score	Token F1 score	Subgroup AUC
BERT (6/12)	0	63.46 $\pm$ 0.21	34.4 $\pm$ 3.87	0.59 $\pm$ 0.01
	0.01	65.12 $\pm$ 0.38	36.3 $\pm$ 4.01	0.707 $\pm$ 0.11
	0.1	65.92 $\pm$ 0.24	39.26 $\pm$ 3.91	0.784 $\pm$ 0.07
DistilBERT (3/6)	0	64.31 $\pm$ 0.83	39.81 $\pm$ 4.22	0.768 $\pm$ 0.24
	0.01	64.35 $\pm$ 0.51	40.4 $\pm$ 3.04	0.748 $\pm$ 0.16
	0.1	65.11 $\pm$ 0.7	41.03 $\pm$ 3.28	0.794 $\pm$ 0.31
RoBERTa (6/12)	0	66.71 $\pm$ 0.22	42.67 $\pm$ 3.14	0.796 $\pm$ 0.28
	0.01	80.86 $\pm$ 0.22	33.19 $\pm$ 3.28	0.612 $\pm$ 0.29
	0.1	78.58 $\pm$ 0.23	36.49 $\pm$ 4.11	0.681 $\pm$ 0.17
DistilRoBERTa (3/6)	0	71.05 $\pm$ 0.43	32.6 $\pm$ 4.01	0.62 $\pm$ 0.08
	0.01	79.14 $\pm$ 0.47	34.41 $\pm$ 4.11	0.634 $\pm$ 0.04
	0.1	81.25 $\pm$ 0.33	36.51 $\pm$ 3.5	0.635 $\pm$ 0.08
	1	81.96 $\pm$ 0.51	43.02 $\pm$ 4.14	0.65 $\pm$ 0.09

$$Loss_{\Sigma} = Loss_{pred} + \lambda Loss_{attn}$$

Performance and fairness scores (Subgroup AUC) of models trained with word-level supervision

BERT Subgroup AUC scores

- .59 - without attention supervision
- .80 - with attention supervision

\* $\lambda = 0$  - non-supervised attention learning

28

## Conclusion on this work

- We conducted two chains of experiments to analyze the effect of Transformer LMs **pruning** in the context of **hate speech classification** tasks (with and without attention supervision)
- We compare **both fairness and performance loss** for pruned BERT, RoBERTa, and their distilled versions
- We show and statistically prove that **removing any layer** from Transformer LMs **results in fairness loss** even when the performance loss could be negligible
- We conduct supervised attention-learning experiments that help to **reduce bias in pruned models**

29

## Recent and ongoing work on LLMs (1)

- Recent work @ERIC:
  - Participation to the **BabyLM challenge** (<https://babylm.github.io>), shared task of CoNLL 2023. Can we pretrain a model from few data? What kind of task can we still solve?  
Our contribution: "Mini Minds: Exploring Bebeshka and Zlata Baby Models"
  - Work on quantization of LLMs (GPTQ) and **confidence**  
Is there an impact of quantization on the model confidence?  
Our contribution: "When Quantization Affects Confidence of Large Language Models?" (NAACL Findings 2024)

30

## Recent and ongoing work on LLMs (2)

- Ongoing work:
  - French dataset related to **ethical questions**, semi-automatic translation of Moral Stories (Emelin et al., 2021), joint work with LabHC
  - French dataset for **implicit hate speech detection**, joint work with Naver Labs
  - new methods for **debiasing** compressed generative (causal) models
  - study the fairness of LMs in **detecting subtle forms** of hate speech (joint work with C. Favre)

31

## How fair are language models in detecting subtle forms of hate speech? (1)

- Create a set of synthetically generated sentences
  - Chat-GPT generated sentences following instructions based on (manual) patterns such as <They> <are good at> <object> and paraphrased in various ways  
e.g.: *The girls are marketing ninjas. Even under pressure, the girls maintain their composure.*
- Evaluate the difference of toxicity scores for various groups of people
  - Difference between the toxicity score for « They » and the plural noun (e.g. « The girls »)
  - Compute the mean for different models fine-tuned on toxicity (e.g., Hate-BERT)

32

## How fair are language models in detecting subtle forms of hate speech? (2)

- Some results:

Attribute	Sensitivity		Toxic Samples	
	<i>BERT</i>	<i>RoBERTa</i>	<i>BERT</i>	<i>RoBERTa</i>
girls	0.093	0.005	44	2
women	0.070	0.008	44	5
boys	0.098	0.006	71	2
men	0.073	0.006	50	4
buddhists	0.112	0.007	20	2
muslims	0.101	0.013	24	9
jews	0.093	0.052	26	30
christians	0.169	0.064	128	32
atheists	0.143	0.027	110	12

33

## Take-away message

- Transformer-based language models **can be used** for detecting hate speech in texts
- language models **can be biased** when detecting hate speech: neutral sentence with phrase « Indian guy » is classified as hateful by LMs
- **lack of datasets in French** for hate speech detection (we are working on it)
- **Compression can amplify bias** in LMs used for hate speech detection (even though biases are already present in pre-compressed models)
- **Bias can be mitigated** in language models with forcing models not to pay attention to gender/nationality/religion expressions.

34



## Some references

- [1] Kenton, Jacob Devlin Ming-Wei Chang, and Lee Kristina Toutanova. "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding." In *Proceedings of NAACL-HLT*, pp. 4171-4186. 2019.
- [2] Hooker, Sara, et al. "What do compressed deep neural networks forget?." *arXiv preprint arXiv:1911.05248* (2019)
- [3] Mathew, Binny, Punyajoy Saha, Seid Muhie Yimam, Chris Biemann, Pawan Goyal, and Animesh Mukherjee. "Hatexplain: A benchmark dataset for explainable hate speech detection." In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 17, pp. 14867-14875. 2021.
- [4] Gupta, Manish, Vasudeva Varma, Sonam Damani, and Kedhar Nath Narahari. "Compression of deep learning models for NLP." In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, pp. 3507-3508. 2020.

35

# Thank you

Email: [Julien.Velcin@univ-lyon2.fr](mailto:Julien.Velcin@univ-lyon2.fr)

Website: <https://eric.univ-lyon2.fr/jvelcin/>

DIKé project website: <https://www.anr-dike.fr>



Projet DIKé