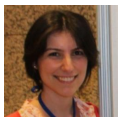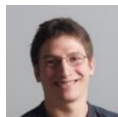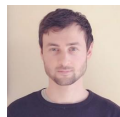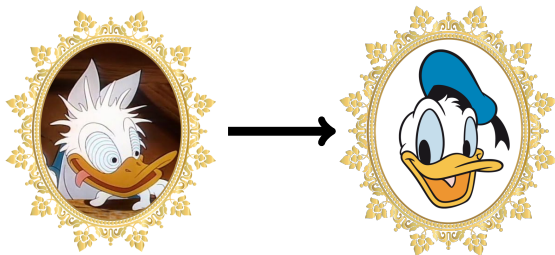# MultiLexNorm: A Shared Task on Multilingual Lexical Normalization

Rob van der Goot, Alan Ramponi, Arkaitz Zubiaga, Barbara Plank,
Benjamin Muller, Iñaki San Vicente Roncal, Nikola Ljubešić,
Özlem Çetinoğlu, Rahmad Mahendra, Talha Çolakoğlu,
Timothy Baldwin, Tommaso Caselli and Wladimir Sidorenko

# Lexical Normalization

Lexical normalization is the task of transforming an utterance into its standard form, word by word, including both one-to-many (1-n) and many-to-one (n-1) replacements.

# Lexical Normalization

State before shared task:

- ▶ Most work on English
- ▶ Also work on single other languages
- ▶ Varieties in task definitions, guidelines and metrics
- ▶ No common evaluation benchmark

# MultiLexNorm

- combination of existing datasets
- annotation style and file format converged
- "new" evaluation metric
- external evaluation (UD)

# MultiLexNorm

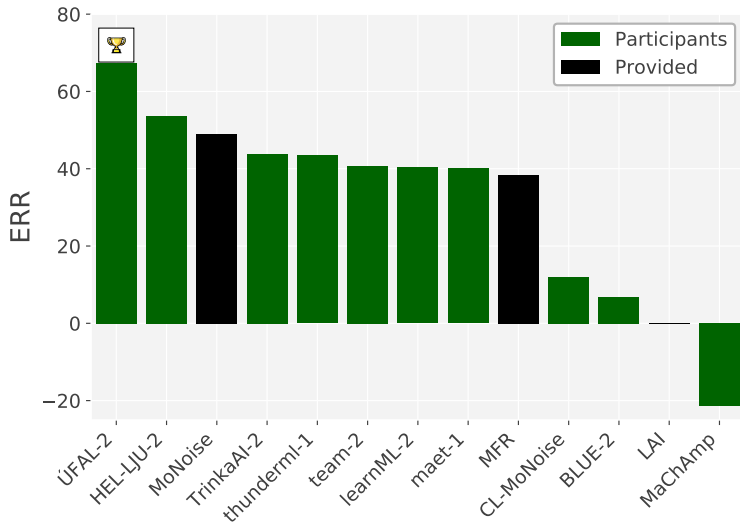| Lang. | Language name | Normalization example |
|-------|---------------|-----------------------|
| DA | Danish | De skarpe lamper gjorde destromindre ek bedre .<br>De skarpe lamper gjorde destro mindre ikke bedre . |
| DE | German | ogäj isch hätts auch dwiddern könn<br>Okay ich hätte es auch twittern können |
| EN | English | u hve to let ppl decide what dey want to do<br>you have to let people decide what they want to do |
| ES | Spanish | @username cuuxamee sii peroo veen yaa eem<br>@username escúchame sí pero ven ya eh |
| HR | Croatian | svi frendovi mi nešto rade , veceras san osta sam .<br>svi frendovi mi nešto rade , večeras sam ostao sam . |
| ID-EN | Indonesian-English | pdhal not fully bcs those ppl jg sih .<br>padahal not fully because those people juga sih . |
| IT | Italian | a Roma è cosí primavera che sembra gia giov<br>a Roma è così primavera che sembra già giovedì |
| NL | Dutch | Kga me wss trg rolle vant lachn<br>Ik ga me waarschijnlijk terug rollen van het lachen |
| SL | Slovenian | jst bi tud najdu kovanec vreden veliko denarja .<br>jaz bi tudi našel kovanec vreden veliko denarja . |
| SR | Serbian | komunalci kace pocne kaznjavanje ?<br>komunalci kad počne kažnjavanje ? |
| TR | Turkish | He o dediyin suala cvb verdim<br>He o dediğin suale cevap verdim |
| TR-DE | Turkish-German | @username Yerimm senii , damkee schatzymm :-*<br>@username Yerim seni , danke Schatzym :-* |

Is this sample biased?

# Metric

- ▶ Previously: accuracy, accuracy over OOV words, F1 score, BLEU, word error rate, character error rate, etc.
- ▶ Now: accuracy normalized for amount of words to be normalized. Error Reduction Rate:
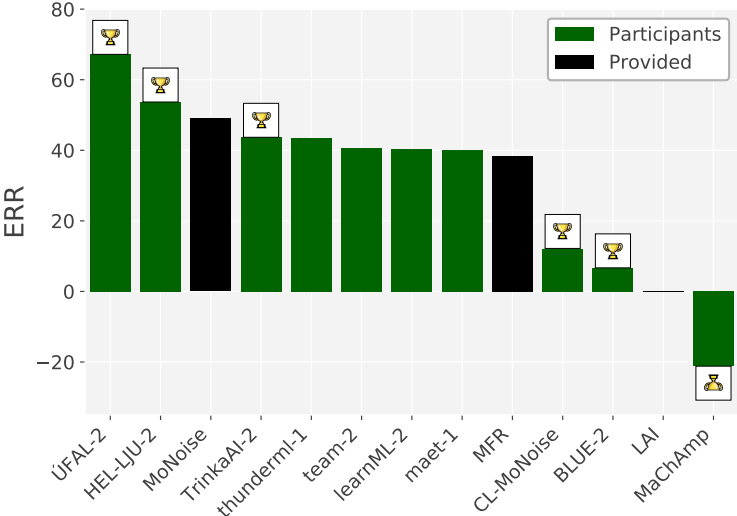
$$ERR = \frac{\%\text{accuracy} - \%\text{words\_not\_normed}}{100 - \%\text{words\_not\_normed}}$$
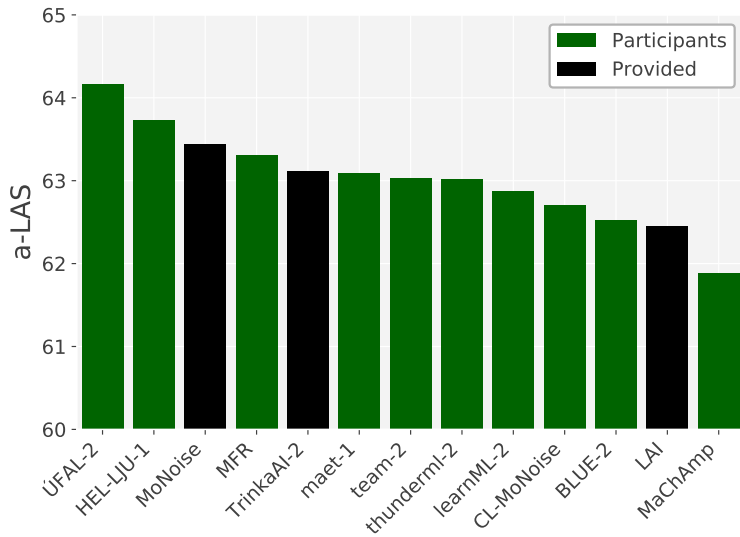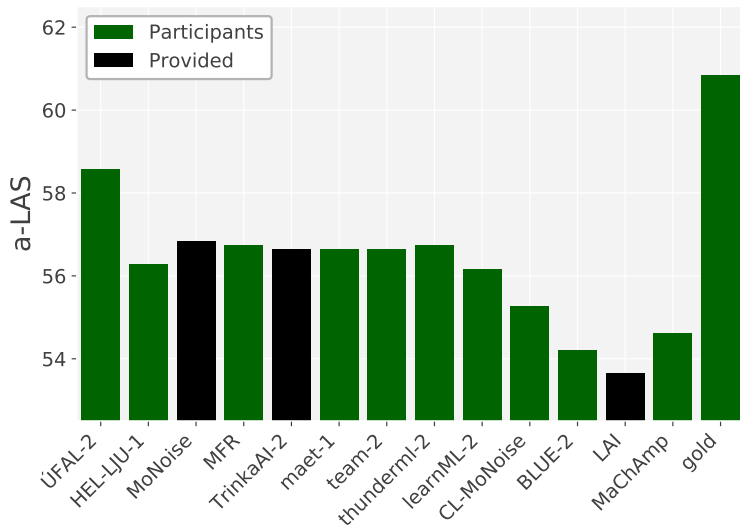
# Results

# Results

# Extrinsic Evaluation (avg.)

# Extrinsic Evaluation (EN-MoNoise)

# Thanks!

- https://bitbucket.org/robvanderg/multilexnorm/