

# Visualizing Neural Machine Translation Attention and Confidence

Matīss Rikters<sup>a</sup>, Mark Fishel<sup>b</sup>, Ondřej Bojar<sup>c</sup>

<sup>a</sup>Faculty of Computing, University of Latvia

<sup>b</sup>Institute of Computer Science, University of Tartu

<sup>c</sup>Charles University, Faculty of Mathematics and Physics, Institute of Formal and Applied Linguistics

E-mail: matiss@lielakeda.lv, fishel@ut.ee, bojar@ufal.mff.cuni.cz

## Abstract

In this article, we describe a tool for visualizing the output and attention weights of neural machine translation systems and for estimating confidence about the output based on the attention.

Our aim is to help researchers and developers better understand the behaviour of their NMT systems without the need for any reference translations. Our tool includes command line and web-based interfaces that allow to systematically evaluate translation outputs from various engines and experiments. We also present a web demo of our tool with examples of good and bad translations: <http://ej.uz/nmt-attention>.

## Confidence Scores

$$CDP = \frac{1}{J} \sum_j \log \left( 1 + \left( \sum_i \alpha_{ji} \right)^2 \right)$$

$$AP_{out} = -\frac{1}{I} \sum_i \sum_j \alpha_{ji} \cdot \log \alpha_{ji}$$

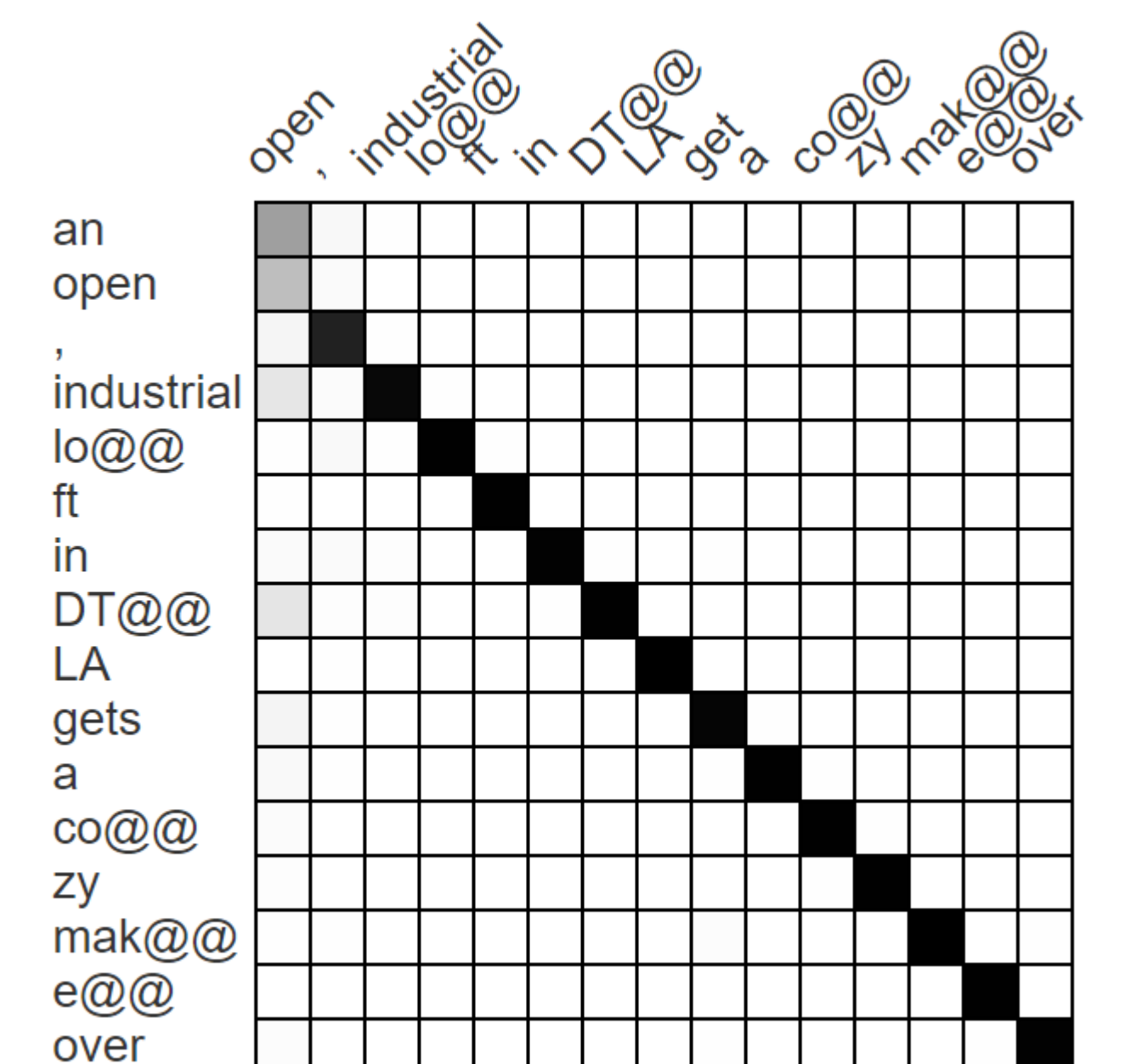
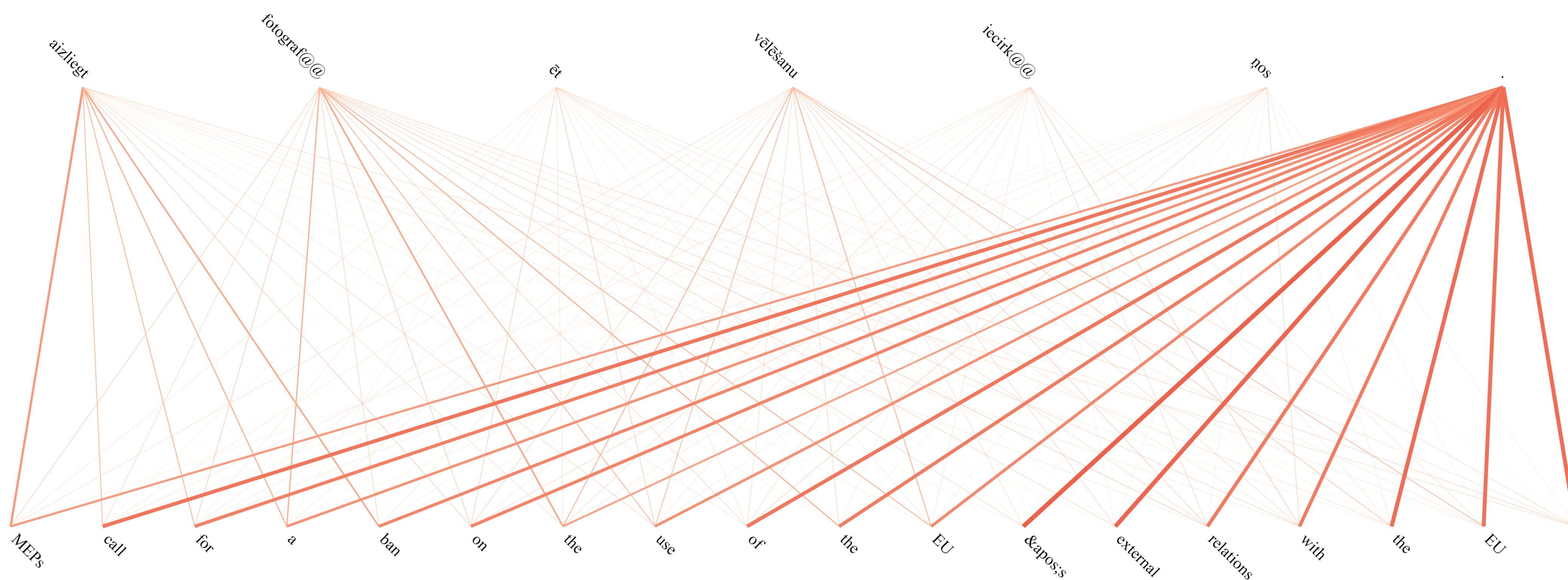
$$confidence = CDP + AP_{out} + AP_{in}$$

$$AP_{in} = -\frac{1}{I} \sum_j \sum_i \alpha_{ij} \cdot \log \alpha_{ij}$$

$$percentage = e^{-C(X^2)}$$

## Lacking Confidence

## Excessive Confidence



## Features

## GitHub Link

- Works with attention alignment data from
  - Nematus
  - Neural Monkey
  - AmuNMT
- Visualise translations in
  - Linux Terminal or Windows PowerShell
  - Web browser
    - Line form or matrix form
    - Save as PNG
    - Sort and navigate dataset by confidence scores



<http://ej.uz/nmt-attention>

## Poster Link



<http://ej.uz/nmt-attention>

## Acknowledgements

P A R S E M E



This research was supported by the ICT COST Action IC1207 ParseME: Parsing and multi-word expressions - towards linguistic precision and computational efficiency in natural language processing, the grant H2020-ICT-2014-1-645442 (QT21) and Charles University Research Programme "Progres" Q18+Q48.