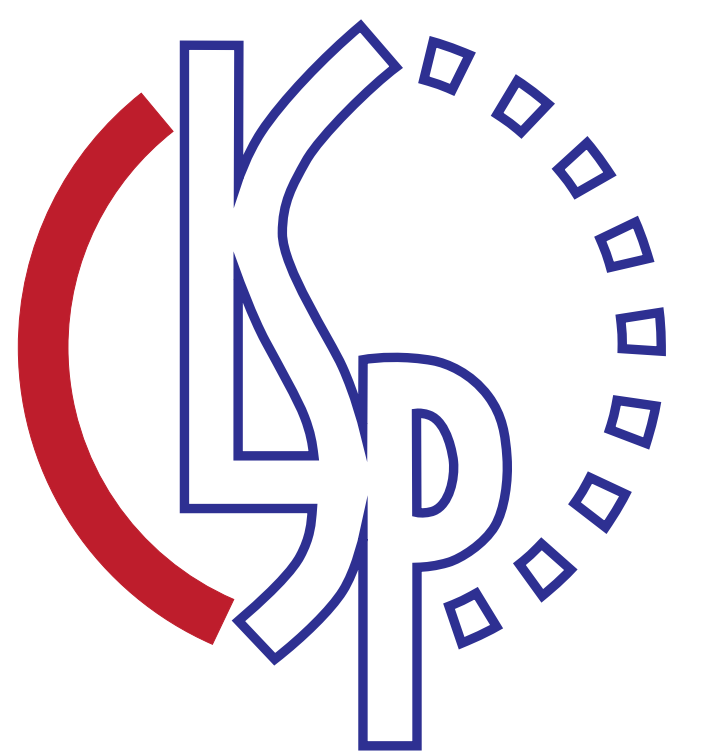# On the Impact of Various Types of Noise on Neural Machine Translation

## Huda Khayrallah & Philipp Koehn
Center for Language & Speech Processing, Computer Science Department
Johns Hopkins University
{huda, phi}@jhu.edu

## Abstract

We examine how various types of noise in the parallel training data impact the quality of neural machine translation systems. We create five types of artificial noise and analyze how they degrade performance in neural and statistical machine translation. We find that neural models are generally more harmed by noise than statistical models. For one especially egregious type of noise they learn to just copy the input sentence.

## Motivation

- MT systems (especially NMT) are data hungry
- Is all data helpful?
- Adding noisy web-crawl hurts **NMT** and helps **SMT**

|             | **NMT**       | **SMT**       |
|-------------|---------------|---------------|
| WMT17       | 27.2          | 24.0          |
| + noisy corpus | 17.3 (-9.9) | 25.2 (+1.2) |

What kind of noise exists in the noisy corpus?

| Type of noise          |      |
|------------------------|------|
| Okay                   | 23%  |
| Misaligned sentences   | 41%  |
| 3rd Language           | 3%   |
| Both English           | 10%  |
| Both German            | 10%  |
| Untranslated           | 4%   |
| Short Segments ≤ 2     | 1%   |
| Short Segments 3 - 5   | 5%   |
| Other Text             | 2%   |

## Experiments

- SMT (Moses) and NMT (Marian)
- WMT 2017 DE → EN baseline
- Web-crawled data from paracrawl.eu

## Results

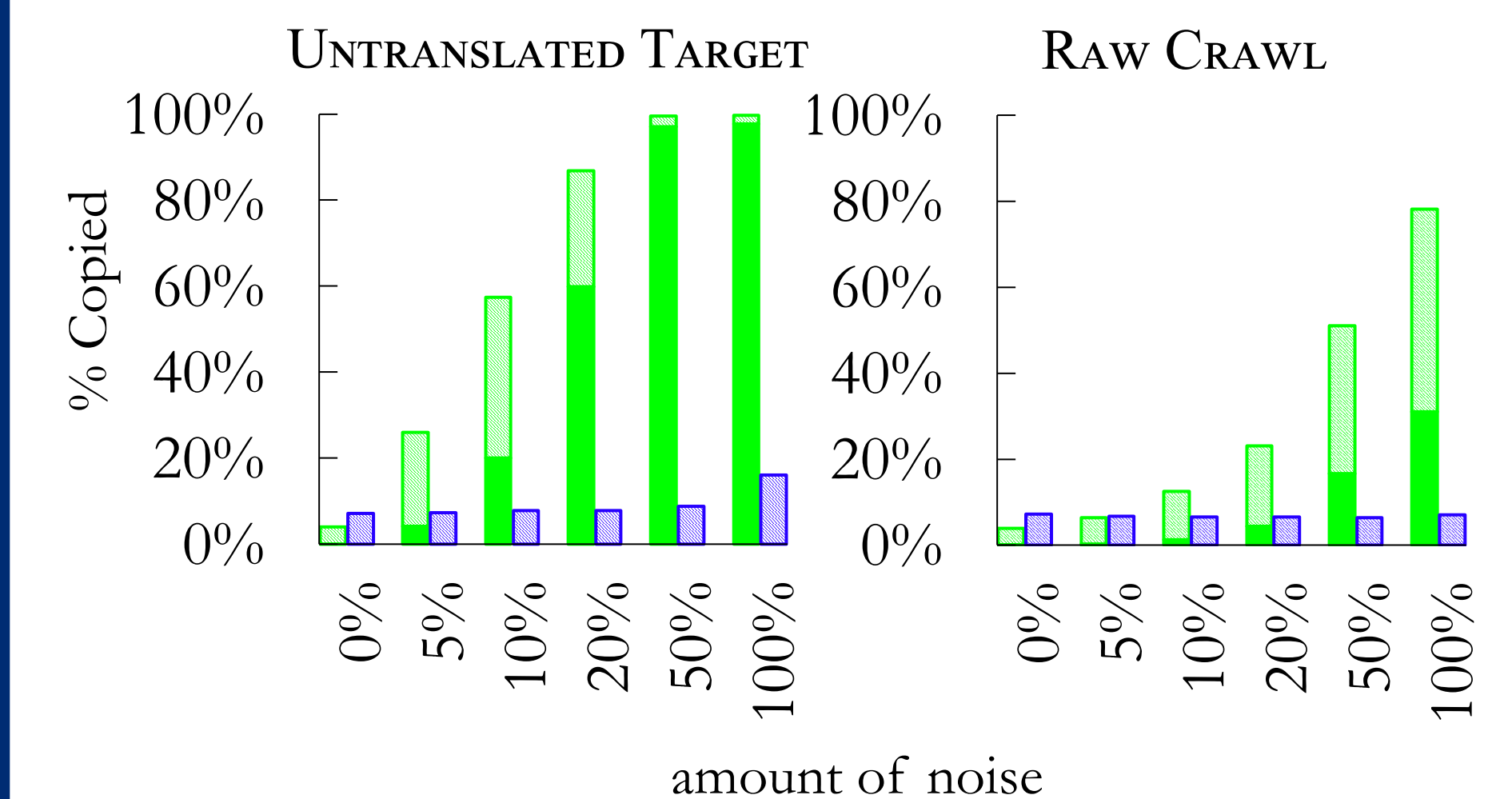BLEU scores after adding different amounts of noise (ratio of original clean corpus)

**NMT (left green bars)** is harmed more than **SMT (right blue bars)**

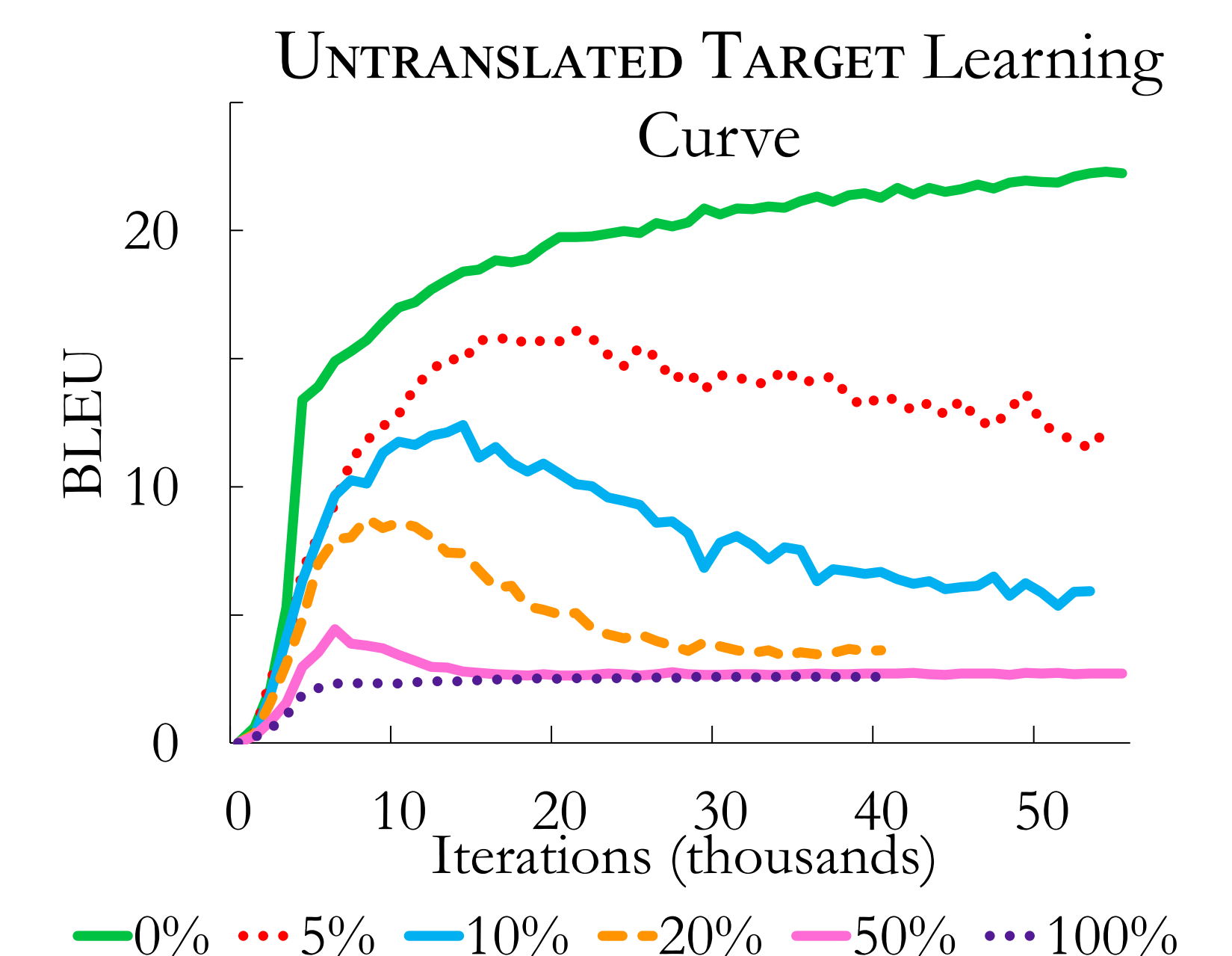|                               | 5%              | 10%             | 20%             | 50%             | 100%            |
|-------------------------------|-----------------|-----------------|-----------------|-----------------|-----------------|
| **MISALIGNED SENTENCES**      | 26.5 24.0 / -0.7 -0.0 | 26.5 24.0 / -0.7 -0.0 | 26.3 23.9 / -0.9 -0.1 | 26.1 23.9 / -1.1 -0.1 | 25.3 23.4 / -1.9 -0.6 |
| **MISORDERED WORDS (SOURCE)** | 26.9 24.0 / -0.3 -0.0 | 26.6 23.6 / -0.6 -0.4 | 26.4 23.9 / -0.8 -0.1 | 26.6 23.6 / -0.6 -0.4 | 25.5 23.7 / -1.7 -0.3 |
| **MISORDERED WORDS (TARGET)** | 27.0 24.0 / -0.2 -0.0 | 26.8 24.0 / -0.4 -0.0 | 26.4 23.4 / -0.8 -0.6 | 26.7 23.2 / -0.5 -0.8 | 26.1 22.9 / -1.1 -1.1 |
| **WRONG LANGUAGE (SOURCE)**   | 26.9 24.0 / -0.3 -0.0 | 26.8 23.9 / -0.4 -0.1 | 26.8 23.9 / -0.4 -0.1 | 26.8 23.9 / -0.4 -0.1 | 26.8 23.8 / -0.4 -0.2 |
| **WRONG LANGUAGE (TARGET)**   | 26.7 24.0 / -0.5 -0.0 | 26.6 23.9 / -0.6 -0.1 | 26.7 23.8 / -0.5 -0.2 | 26.2 23.5 / -1.0 -0.5 | 25.0 23.4 / -2.2 -0.6 |
| **UNTRANSLATED (SOURCE)**     | 27.2 23.9 / -0.0 -0.1 | 27.0 23.9 / -0.2 -0.1 | 26.7 23.6 / -0.5 -0.4 | 26.8 23.7 / -0.4 -0.3 | 26.9 23.5 / -0.3 -0.5 |
| **UNTRANSLATED (TARGET)**     | 17.6 23.8 / -9.8 -0.2 | 11.2 23.9 / -16.0 -0.1 | 5.6 23.8 / -21.6 -0.2 | 3.2 23.4 / -24.0 -0.6 | 3.2 21.1 / -24.0 -2.9 |
| **SHORT SEGMENTS (max 2)**    | 27.1 24.1 / -0.1 +0.1 | 26.5 23.9 / -0.7 -0.1 | 26.7 23.8 / -0.5 -0.2 |                 |                 |
| **SHORT SEGMENTS (max 5)**    | 27.8 24.2 / +0.6 +0.2 | 27.6 24.5 / +0.4 +0.5 | 28.0 24.5 / +0.8 +0.5 | 26.6 24.2 / -0.6 +0.2 |                 |
| **RAW CRAWL DATA**            | 27.4 24.2 / +0.2 +0.2 | 26.6 24.2 / -0.6 +0.2 | 24.7 24.4 / -2.5 +0.4 | 20.9 24.8 / -6.3 +0.8 | 17.3 25.2 / -9.9 +1.2 |

## Analysis

What goes wrong with UNTRANSLATED TARGET?

- NMT learns to copy input, causing degradation.



- Copied sentences in the UNTRANSLATED TARGET & RAW CRAWL experiments for **NMT (left green bars)**, and **SMT (right blue bars)**.

- Sentences that are exact matches to the source are the solid bars, sentences that are more similar to the source than the target are the shaded bars.



UNTRANSLATED TARGET Learning Curve

Legend: 0%, 5%, 10%, 20%, 50%, 100%

- Performance of the systems trained on noisy corpora begin to improve, before over-fitting to the copy portion of the training set.