

Sequence Length is a Domain: Length-based Overfitting in Transformer Models

Dušan Variš and Ondřej Bojar

Charles University, Institute of Formal and Applied Linguistics

October 11th, 2021

Overfitting in Transformers

- ▶ Recent models (e.g. GPT-3) increase both in size and in number of training instances.
- ▶ We suspect that an overlap in the train-test data could lead to overestimation of model generalization ability.
- ▶ Long-range dependencies in transformer:
 - ▶ result of poor modeling ability (?) ...
 - ▶ ... or lack of data with long-range dependencies?

Mock Task: String Editing

- ▶ Easier evaluation:
 - ▶ clear distinction between examples,
 - ▶ no ambiguity in correct answers,
 - ▶ accuracy metric: exact match with correct solution

Input	Output
push 1 1 0 1 0	1 0 1 0 1
reverse - 1 0 0 1 1	1 1 0 0 1

String Editing: Results

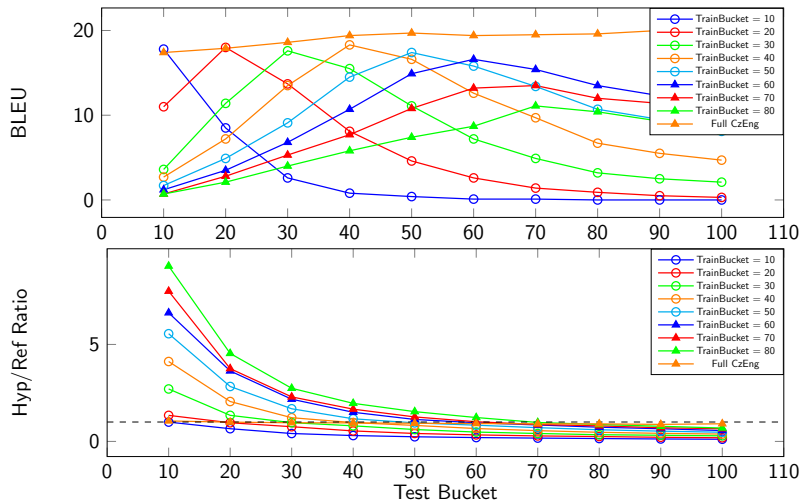
- ▶ Training lengths: 10-15
- ▶ train/test sequence length mismatch → the models fail horribly

	0-10	11-15	16-20
copy	62.6	100.0	0.0
push	59.1	100.0	0.0
pop	0.1	100.0	0.0
shift	52.5	100.0	0.0
unshift	41.2	100.0	0.0
reverse	0.0	84.4	0.0
all	15.822	97.5	0.978

Machine Translation

- ▶ Split CzEng 2.0 (Kocmi et al., 2020) into buckets based on target-side (or source-side) sequence length (after subword tokenization).
- ▶ Train a separate system on each training bucket.
- ▶ Evaluate on WMT newstest split in a similar way.

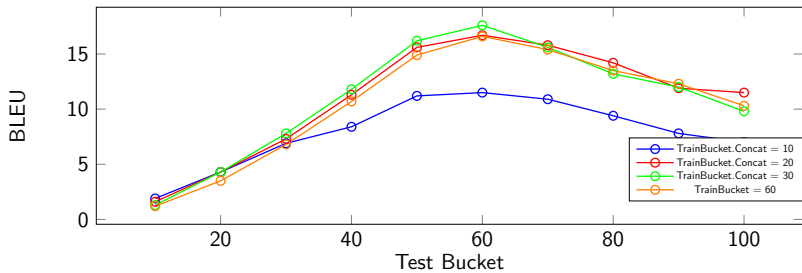
Machine Translation: Results (Target-length Buckets)



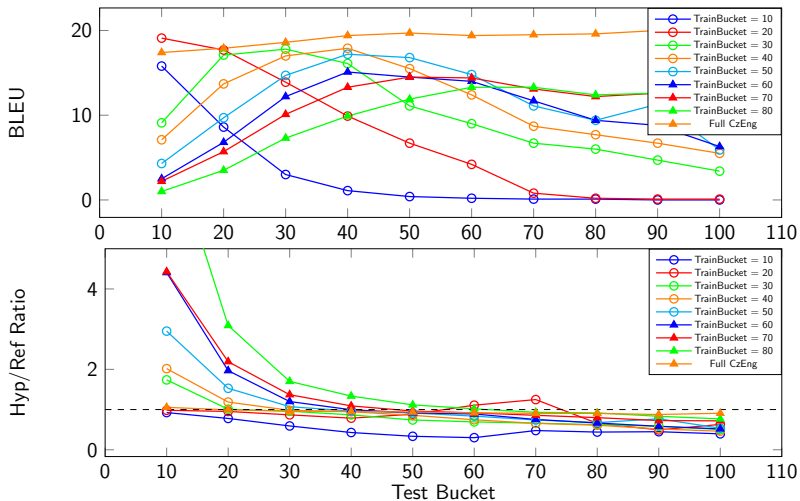
Machine Translation: Synthetic concatenated data

- ▶ Create a synthetic 60-bucket data using concatenation of:
 - ▶ 6 x 10-bucket sentences,
 - ▶ 3 x 20-bucket sentences,
 - ▶ 2 x 30-bucket sentences.
- ▶ We concatenate consecutive sentence pairs (after shuffling).
- ▶ Compare with the system trained on 60-bucket data.

Synthetic Concatenation: Results (Target-length Buckets)



Machine Translation: Source-length Buckets



Source-length Buckets: Target-length Distributions

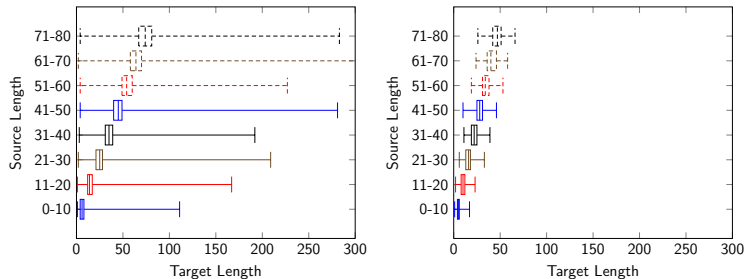


Figure: **Left:** Train Distribution, **Right:** Test Distribution

See You at the Posters!

Sequence Length is a Domain: Length-based Overfitting in Transformer Models

Dušan Varšić and Ondřej Bojar
varis.bojar@ufal.mff.cuni.cz

Charles University, Faculty of Mathematics and Physics
Institute of Formal and Applied Linguistics
Prague, Czech Republic



Introduction

Transformers generalize poorly to longer AND shorter sequence editing examples. Similar trends can be observed on MT task.

Methods

Split data to buckets based on target-side length. Train a separate NMT system on each training bucket and evaluate it on the validation buckets.

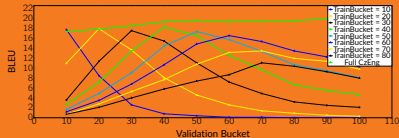
Results

Strong implication of target-side-length overfitting in Transformers that use absolute position encoding.

- Higher train-test length difference → higher performance drop.
- Hypothesis length similar to that of training data.
- Length overfitting could be avoided with relative position embeddings (Neishi and Yoshinaga, 2019).

Transformers with absolute position encoding output sequences of length similar to sequences in training data.

Source (30-bucket)	The company does not collect its mail and it has closed its official headquarters in Žilkov more than six years ago.
Hyp1 (10-bucket)	Společnost nebere poštu a zavřela úřední sídlo.
Hyp1 (30-bucket)	The company does not gather mail and closed official headquarters.
Hyp2 (30-bucket)	Společnost nebere poštu a již před více než šesti lety zavřela své oficiální sídlo v Žilkově.
Hyp2 (30-bucket)	The company does not collect mail and more than six years ago closed its official headquarters in Žilkov.
Hyp3 (60-bucket)	Společnost nebere poštu a uzavřela své oficiální sídlo v Žilkově více než šest let ago. v Žilkově. Společnost nebere poštu a uzavřela oficiální sídlo v Žilkově více než šest let ago. v Žilkově.
Hyp3 (30-bucket)	The company does not gather up mail and closed up its official its official headquarters in Žilkov more than six years ago. in Žilkov. The company does not collect mail and closes up official headquarters in Žilkov more than six years ago. o.
Reference (30-bucket)	Nebere poštu a oficiální sídlo na Žilkově uzavřel více než šesti lety.
Ref (30-bucket)	(The company) does not collect mail and official headquarters in Žilkov closed up more than six years ago.



This work was supported by the GA ČR MHD14/1: grant Project Representations in Multi-modal and Multi-Signal Modelling, 19-24724X-BYV, GA19-24724X, by the IS4AC grant, Grant Science & Humanities Open Cloud, 822780 and by SVV 240 452 grant. This template is based on Mike Harrison's idea of Better Poster. It was modified to latex by Rafal Sliwa and modified for UFAL purposes by Tom Kocourek.