

# Improving Fuzzy Match Augmented Neural Machine Translation through Synthetic Data

**Authors:** Arda Tezcan, Alina Skidanova, Thomas Moerman

Tezcan, A., Skidanova, A., & Moerman, T. (2024). Improving Fuzzy Match Augmented Neural Machine Translation in Specialised Domains through Synthetic Data. *The Prague Bulletin of Mathematical Linguistics*, 122, 9–42.

1. Introduction & Motivation
2. Background: Fuzzy Matches, Back-Translation & NMT
3. Proposed Method: Combining Back-Translation & Neural Fuzzy Repair
4. Experimental Setup
5. Results & Analysis
6. Conclusion & Future Work

# Intro & Motivation

# Intro & Motivation

- ◉ **While powerful, NMT (and LLMs) can struggle with specialized domains (legal, technical, ...) due to limited in-domain parallel data** (Jiao et al., 2023)
- ◉ **LLMs are increasingly used for MT, thanks for additional capabilities over NMT models:**
  - Instruction following
  - In-context learning (i.e. few-shot learning)
- ◉ **However, LLMs do not guarantee superior performance over specialized NMT models in domain-specific scenarios** (Wassie et al., 2025; Higashiyama, 2024)

# Intro & Motivation

## ○ **Fuzzy Match (FM) Augmentation:**

- Leveraging translations of similar sentences (FMs) from translation memories or training data improves quality (Bulté & Tezcan, 2019; Xu et al., 2020)
- Challenge: FM effectiveness often depends on *large* bilingual datasets to find *high-similarity* FMs

## ○ **This Study's Goal:**

- Can we improve FM-augmented NMT when in-domain parallel data is *limited*?
- Hypothesis: Yes, by leveraging *additional monolingual data* to generate *synthetic parallel data* specifically for FM augmentation in specialized domains

# Background

# Background: Back-Translation

## **Data Sources in Specialized Domains**



**Limited Parallel Data**

**300K pairs**

High-quality human translations  
In-domain (Legal/Technical)

**HIGH QUALITY**



**Abundant Monolingual**

**1.5M sentences**

Target language only  
In-domain (Legal/Technical)

**HIGH QUALITY**

# Background: Back-Translation

## Data Sources in Specialized Domains



Limited Parallel Data

**300K pairs**

High-quality human translations  
In-domain (Legal/Technical)

HIGH QUALITY



Abundant Monolingual

**1.5M sentences**

Target language only  
In-domain (Legal/Technical)

HIGH QUALITY

## Back-Translation Process

1

**Start with**

High-Quality  
Parallel Data

300K pairs



2

**Train**

Reverse NMT Model  
(Target → Source)

Using 300K pairs



3

**Translate**

Monolingual Data  
Generate Synthetic  
Source

1.5M sentences



4

**Result**

Synthetic Parallel  
Data Pairs

1.5M pairs

# Background: Back-Translation

## Step 1: Original Parallel Data

EN: "The contract is valid"

NL: "Het contract is geldig"

**300K such pairs**

# Background: Back-Translation

**Step 2: Train Reverse Model (NL→EN)**



**Reverse NMT Model**  
**NL → EN**

**Trained on 300K parallel pairs**

# Background: Back-Translation

## Step 3: Apply to Monolingual Dutch Data

### Monolingual Dutch (1.5M sentences)

"Het document moet worden goedgekeurd"

"Deze clausule is belangrijk"

"De voorwaarden zijn onderhandelbaar"

*... and 1.5M more sentences*

NL→EN Model



### Generated English (Synthetic)

"The document must be approved"

"This clause is important"

"The terms are negotiable"

*... 1.5M synthetic sentences*

# Background: Back-Translation

## Step 4: Resulting Synthetic Parallel Data

### New Synthetic Training Pairs

EN: "The document must be approved" (synthetic)

NL: "Het document moet worden goedgekeurd" (original)

EN: "This clause is important" (synthetic)

NL: "Deze clausule is belangrijk" (original)

**1.5M new synthetic parallel pairs**

# Background: Back-Translation

**300K**

Parallel Data

High Quality



**1.5M**

Back-Translated

Synthetic/Noisy

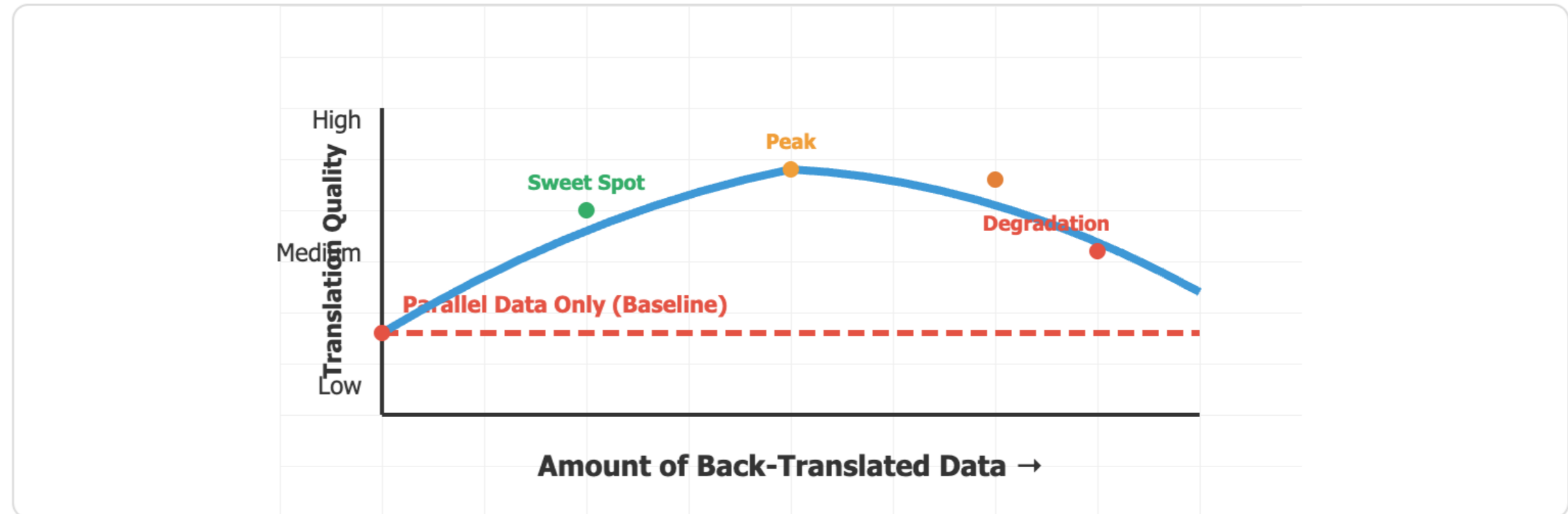


**1.8M**

Total Training

Mixed Quality

## Performance vs. Back-Translation Data Amount



### Optimal Range:

- Back-translation significantly improves translation quality

### Degradation Point:

- Too much synthetic data hurts performance due to noise

# Background: Fuzzy Matches

## What are Fuzzy Matches?

**Fuzzy Matches (FMs)** are similar translations retrieved from translation memory or training data that can guide machine translation systems towards better outputs.

### Example: Finding Fuzzy Matches

**New Source:**

"The contract is invalid"

**85%  
similar**

**Fuzzy Match:**

"Het contract is ongeldig"

**New Source:**

"The document must be signed"

**78%  
similar**

**Fuzzy Match:**

"Het document moet worden ondertekend"

# Background: Fuzzy Matches



## Architecture Modification

Modify the MT model architecture with extra components like attention layers, special decoders, or FM-specific modules.

*Cao and Xiong, 2018; He et al., 2020*



## Shared Goal

**Steer MT output towards translations of fuzzy matches**

All three approaches aim to leverage similar existing translations to improve the quality and consistency of machine translation systems.

*Moslem et al., 2023*

S

## Data Augmentation (IFR)

Repair: Augment training data by inserting source sentences with fuzzy matches. No architecture modification needed.

`<sep> fuzzy_target`

*Reizcan, 2019*

### Step 1: Original Sentence Pair

**Source**

**(S):** We found three studies for inclusion in the review.

**Target**

**(T):** We vonden drie studies voor opname in de review.

### Step 2: Retrieved Fuzzy Match (Similarity: 0.9309)

**FMS:** We identified **nine** eligible studies for inclusion in the review.

**FMT:** We identificeerden **negen geschikte** studies voor opname in de review.

*Differences highlighted in bold*

### Step 3: Creating Augmented Input (S\*)

**Formula:** S\* = Source <sep> FMT

**Actual S\*:**

We found three studies for inclusion in the review.

<sep> We identificeerden negen geschikte studies voor opname in de review.

**Same process applied at inference time - source is augmented with FMT before translation**

Fuzzy Matches Retrieved From



Training Data

OR



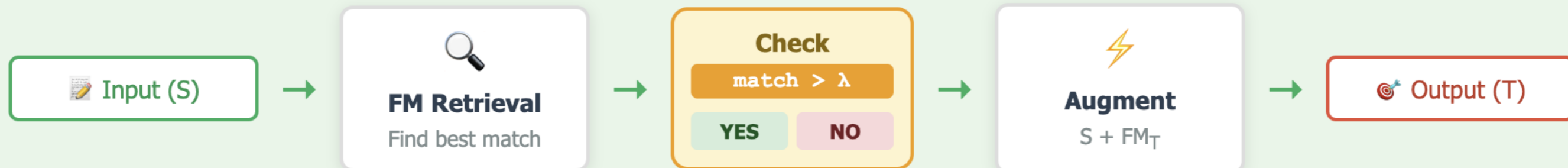
Translation Memory

### 🦾 Training Phase



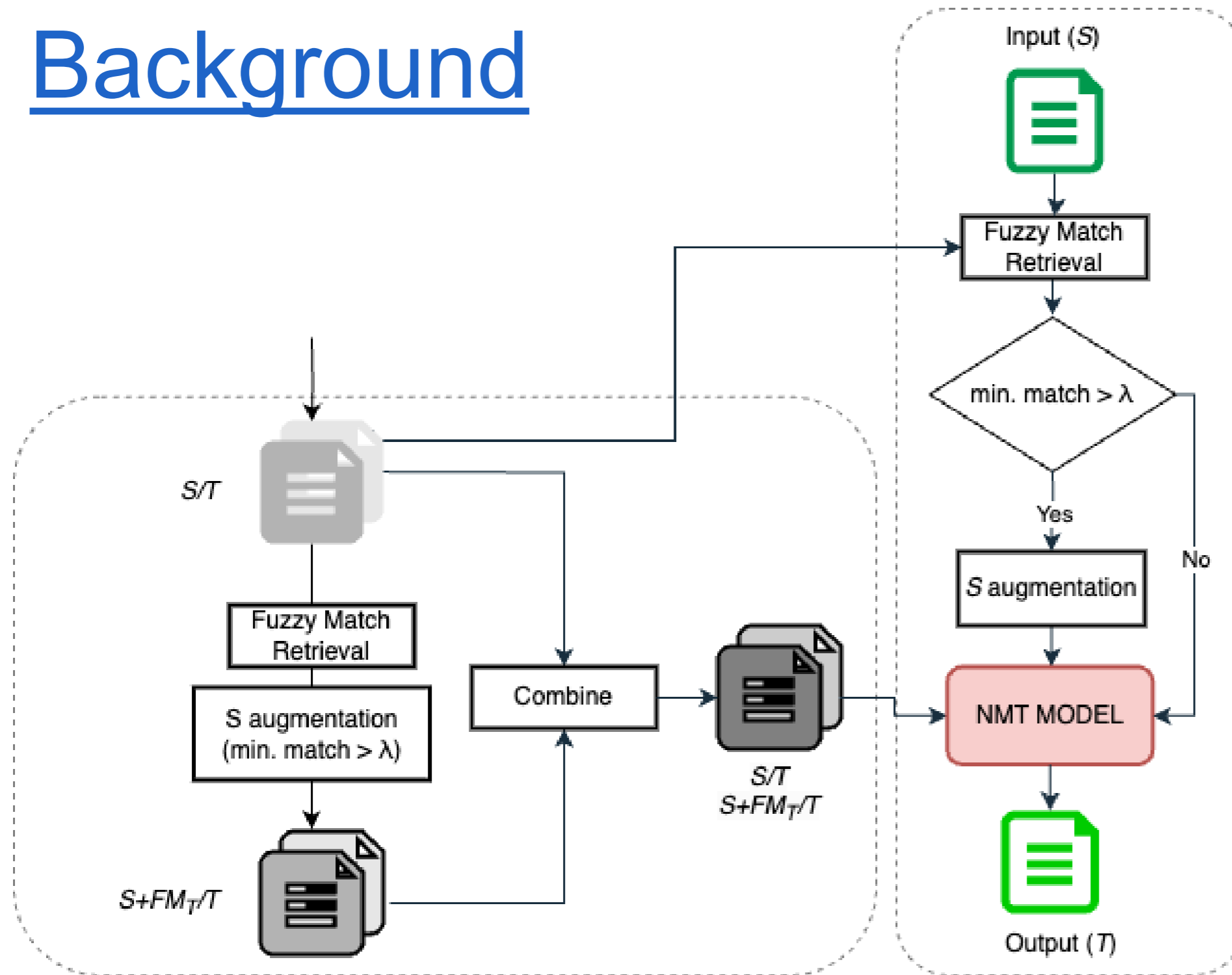
## 🤖 NMT MODEL

### 🎯 Inference Phase



# Background

Towards a better integration of fuzzy matches in neural machine translation through data augmentation  
[Bulté & Tezcan, 2019]



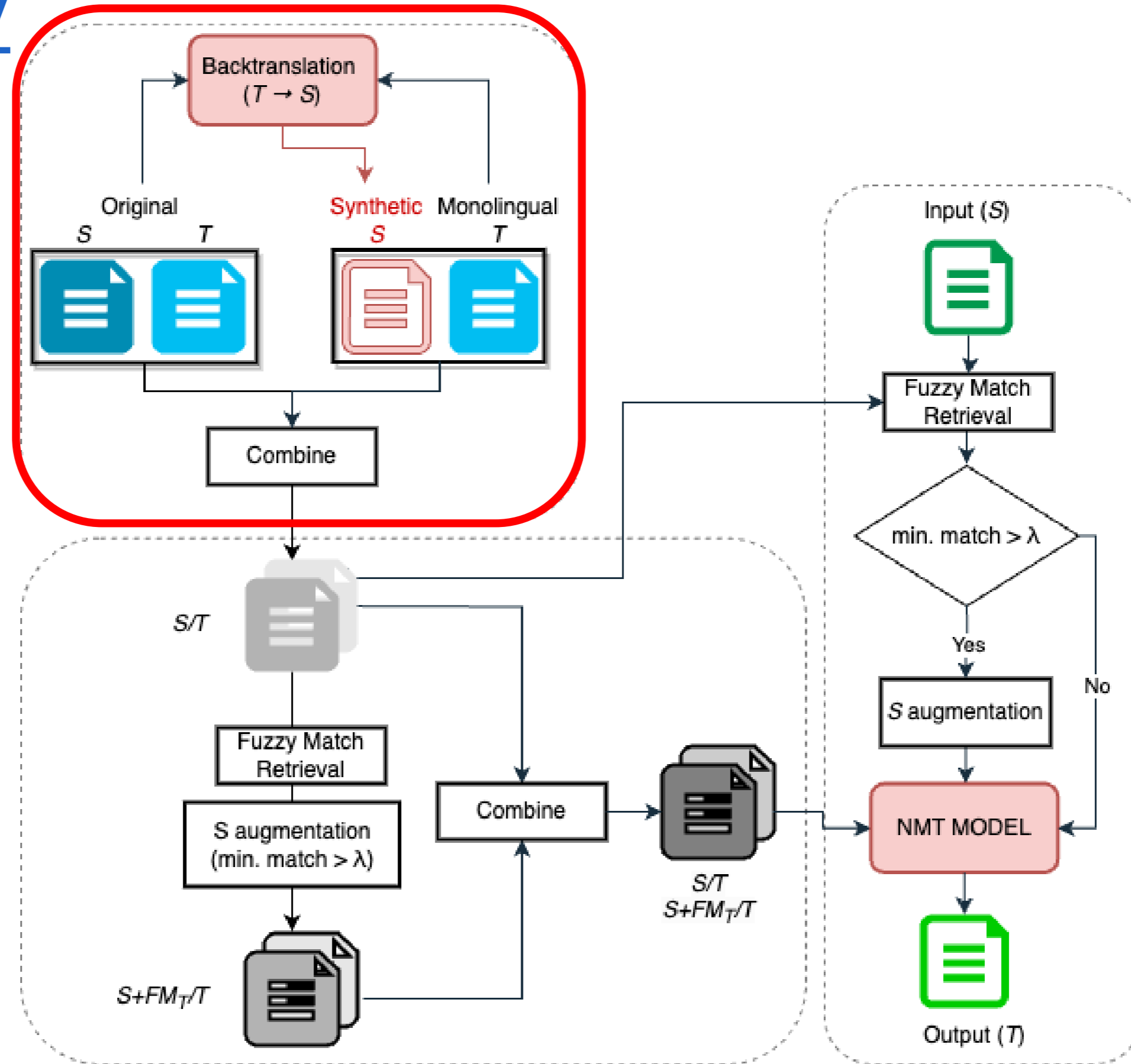
$S$	Debt, breakdown by residual maturity
$score$	0.9812
$FM_S$	Debt, breakdown by <b>initial</b> maturity
$FM_T$	Dette, ventilation par échéance <b>initiale</b>
$S'$	Debt, breakdown by residual maturity < <i>sep</i> > Dette, ventilation par échéance initiale
$T$	Dette, ventilation par échéance résiduelle

# Background

- FM augmentation results in optimal NMT performance in **high-resource scenarios** = (very) large bilingual training data [Bulté & Tezcan, 2019; Xu et al., 2023; Reheman et al., 2023]
- To remedy this limitation, **additional monolingual data in the target language** was used to retrieve FMs directly from texts in target language (via multilingual sentence embeddings) [Tamura et al., 2023]
- Our goal is similar:
  - Extend the effectiveness of FM-augmentation **to lower-resource scenarios**
  - Focus = limited bilingual data, and **additional monolingual data**
- Our approach is different:
  1. MT additional monolingual data to source language (i.e. **back-translation**)
  2. Use the **synthetic bilingual data** as:
    - additional training data
    - for FM augmentation

# Proposed Method: Combining Back-Translation & Neural Fuzzy Repair

# Methodology



## B BASE

BASELINE

Standard NMT training using only original bilingual data.

### Training Data:

Original Bilingual Data ~300K pairs

## BT BT

BASELINE

Back-translation approach: Original data plus synthetic bilingual data for training.

### Training Data:

Original Bilingual Data ~300K pairs

Back-translated Data ~1.5M pairs

### Total Training:

~1.8M sentence pairs

## BT+NFI BT + NFR

PROPOSED

**Our Approach:** Original + additional synthetic bilingual data used for both training AND fuzzy match augmentation.

### Training Data:

Original Bilingual Data ~300K pairs

Back-translated Data ~1.5M pairs

+ FM Augmentation

### Key Innovation:

Synthetic data for Training + FM Pool

Total: ~1.8M pairs with FM augmentation

## NFR NFR

BASELINE

Neural Fuzzy Repair: Original data with fuzzy match augmentation.

### Training Data:

Original Bilingual Data ~300K pairs

+ FM Augmentation

### Augmentation:

Training + Inference

## NFR NFRmono

BASELINE

NFR with additional monolingual data for FM retrieval during inference only.

### Training Data:

Original Bilingual Data ~300K pairs

### Inference Only:

+ Monolingual Target Data

+ FM Augmentation

### Augmentation:

Inference Only

# Methodology

## Feature Comparison

System	Original Data	Synthetic/HQ Data	FM Augmentation	Key Feature
<b>BASE</b>	✓	✗	✗	Standard NMT baseline
<b>BT</b>	✓	✓ (Synthetic)	✗	Back-translation only
<b>NFR</b>	✓	✗	✓	FM augmentation only
<b>NFRmono</b>	✓	✗	✓ (inference)	FM with monolingual data
<b>BT+NFR (Ours)</b>	✓	✓ (Synthetic)	✓	<b>Combined approach</b>
<b>BASE_HQ</b>	✓	✓ (High-Quality)	✗	<b>Oracle baseline</b>
<b>NFR_HQ</b>	✓	✓ (High-Quality)	✓	<b>Oracle upper bound</b>

# Methodology

## **Training Framework**

**OpenNMT:** Train all NMT models from scratch (no LLMs)

## **Evaluation Metrics**

### **BLEU**

Precision-based metric measuring n-gram overlap between hypothesis and reference translations

### **chrF**

Character-level F-score that considers both precision and recall of character n-grams

### **COMET**

Neural metric trained to predict human judgments of translation quality

## **Statistical Significance Testing**

Bootstrap resampling to ensure reliable performance comparisons

# Data

# DATA

## English → Ukrainian

Legal Domain, Parliament Proceedings

### Sources:

 EU acts • Ukrainian law

 Ukrainian law docs • Legal corpus

	Train	Val	Test
Bilingual	286K	2K	2K
Monolingual	1.46M	—	—

## English ↔ French

European Legislation

### Sources:

 DGT-TM (Bilingual)

 DGT-TM (Monolingual FR + EN)

	Train	Val	Test
Bilingual	300K	2K	2K
Mono (FR)	1.5M	—	—
Mono (EN)	1.5M	—	—

### NFR Limitation

~300K pairs

NFR not useful in low-resource setting

[Bulté & Tezcan, 2019]

### BT Benefits

~1.5M sentences

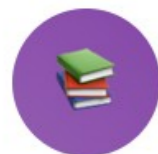
Up to 5x bilingual data improves results

[Edunov et al., 2018]

### Lower-Resource

Test **reduced data sizes** to simulate even more constrained scenarios

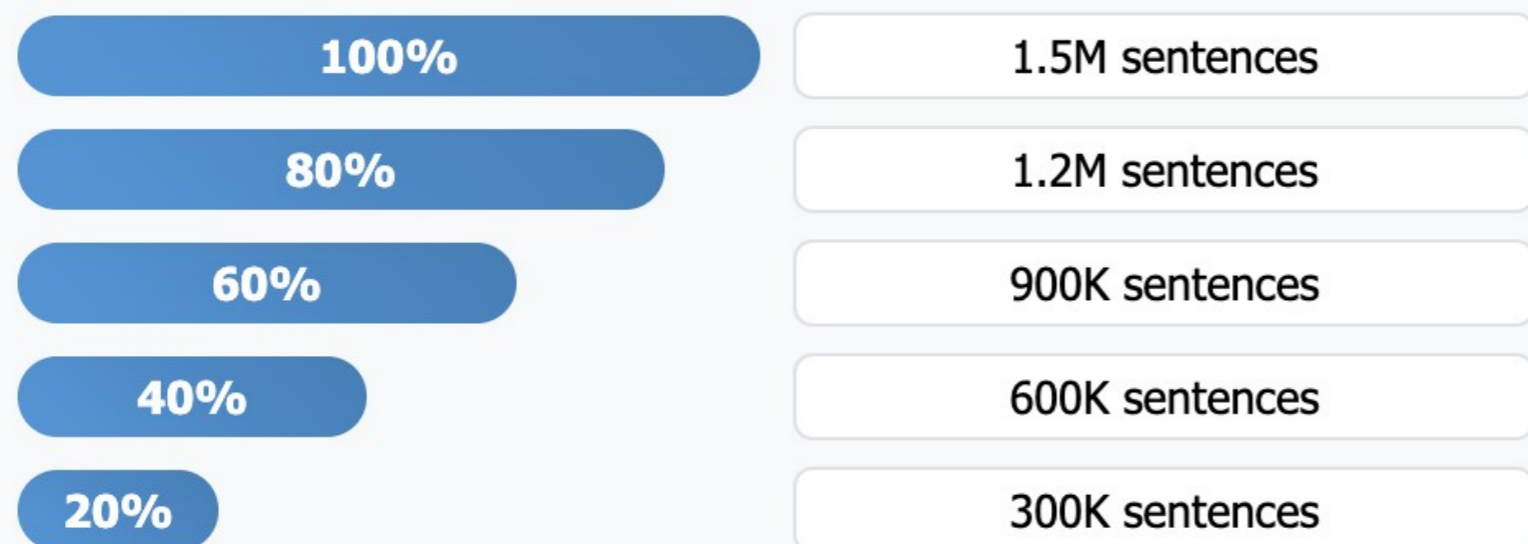
# DATA



## Monolingual Data Reduction

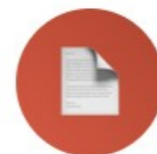
Vary the amount of monolingual data while keeping bilingual data constant at ~300K pairs.

### Monolingual Data Amounts:



### Purpose:

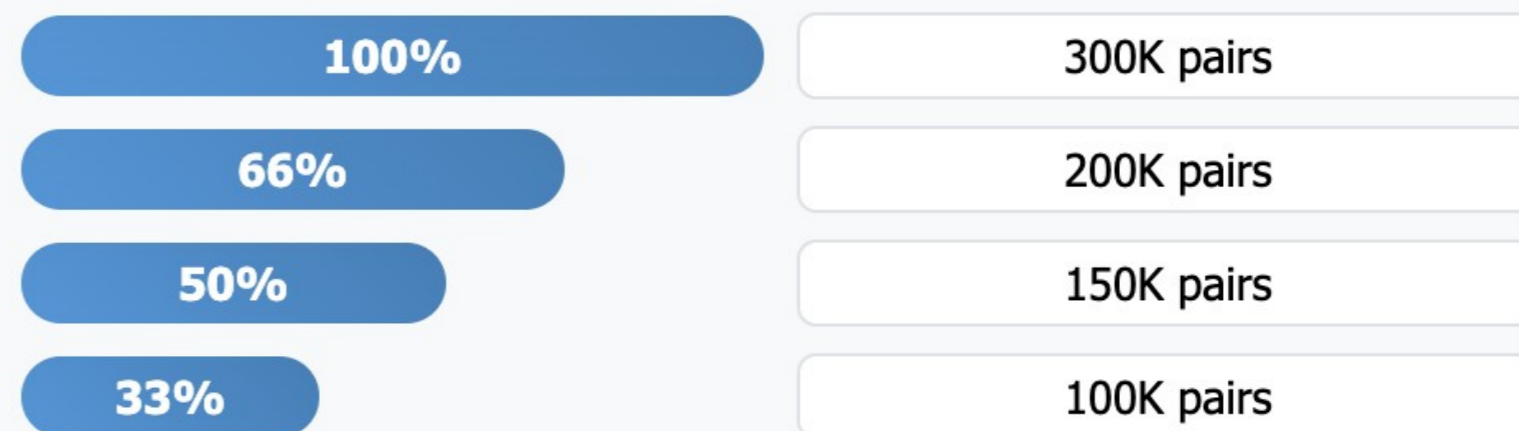
Test how much monolingual data is needed for effective back-translation in FM augmentation.



## Bilingual Data Reduction

Vary the amount of bilingual data while keeping monolingual data constant at ~1.5M sentences.

### Bilingual Data Amounts:



### Purpose:

Test robustness in extremely low-resource scenarios where even less bilingual data is available.

# Results & Conclusion

---

# Results

## 1. Using all available data (bilingual & monolingual)

System	EN→UK			EN→FR			FR→EN		
	BLEU	chrF	COMET	BLEU	chrF	COMET	BLEU	chrF	COMET
<i>BASE</i>	54.90	75.68	90.79	45.76	64.97	79.75	47.76	65.19	80.69
<i>BT (Sennrich et al., 2016)</i>	56.18	76.63	91.61	49.12	67.19	81.42	50.42	67.68	82.32
<i>NFR (Tezcan et al., 2021)</i>	57.71	77.23	91.34	45.61	64.91	79.90	48.05	65.44	80.82
<i>NFR<sub>mono</sub> (Tamura et al., 2023)</i>	60.26	78.58	91.55	46.73	65.48	79.88	49.02	66.05	81.05
<i>Proposed system</i>	<b>66.79</b>	<b>81.99</b>	<b>92.23</b>	<b>57.07</b>	<b>71.83</b>	<b>82.82</b>	<b>60.04</b>	<b>72.96</b>	<b>83.55</b>



### Baseline Analysis

- ✓ **BT** is the best baseline for EN ↔ FR
- ✓ **NFRmono** is the best baseline for EN → UK
- ✓ **NFR** not useful for EN ↔ FR (confirmed)
- ✓ **NFR** useful for EN → UK



### Our System Performance

- ✓ Outperforms **all baseline systems**
- ✓ Up to **+9.62 BLEU** over best baseline
- ✓ Up to **+12.28 BLEU** over BASE
- ✓ Consistent improvements across all language pairs

# Results

## 1. Using all available data (bilingual & monolingual)

System	EN→UK			EN→FR			FR→EN		
	BLEU	chrF	COMET	BLEU	chrF	COMET	BLEU	chrF	COMET
<i>BASE</i>	54.90	75.68	90.79	45.76	64.97	79.75	47.76	65.19	80.69
<i>BT (Sennrich et al., 2016)</i>	56.18	76.63	91.61	49.12	67.19	81.42	50.42	67.68	82.32
<i>NFR (Tezcan et al., 2021)</i>	57.71	77.23	91.34	45.61	64.91	79.90	48.05	65.44	80.82
<i>NFR<sub>mono</sub> (Tamura et al., 2023)</i>	60.26	78.58	91.55	46.73	65.48	79.88	49.02	66.05	81.05
<i>Proposed system</i>	<b>66.79</b>	<b>81.99</b>	<b>92.23</b>	<b>57.07</b>	<b>71.83</b>	<b>82.82</b>	<b>60.04</b>	<b>72.96</b>	<b>83.55</b>
<i>BASE_HQ</i>	–	–	–	53.94	70.27	83.18	54.90	70.38	83.87
<i>NFR_HQ</i>	–	–	–	57.75	72.57	83.91	60.63	73.68	84.48

### NFR\_HQ Performance

**NFR\_HQ performs best** (as expected for high-resource scenarios)

Our proposed system performs **very close** to NFR\_HQ:

- EN→FR: 57.07 vs 57.75 (-0.68 BLEU)
- FR→EN: 60.04 vs 60.63 (-0.59 BLEU)

Using synthetic data instead of high-quality human translations does not drastically impair translation quality.

### BASE\_HQ Comparison

Our proposed system **outperforms BASE\_HQ** in BLEU/chrF:

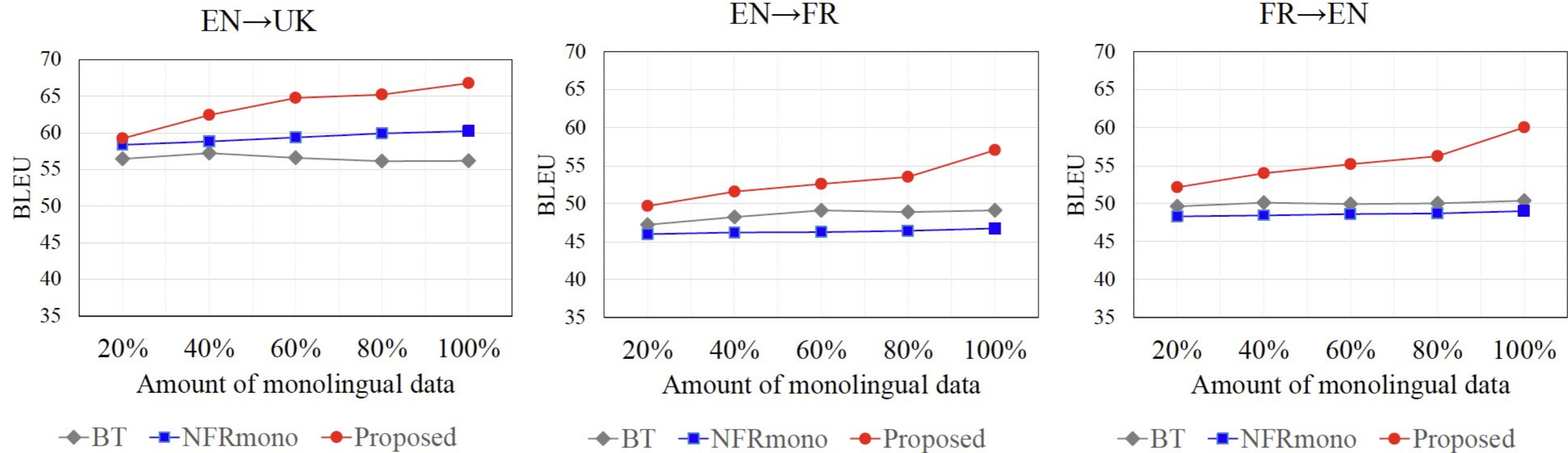
- EN→FR: 57.07 vs 53.94 (+3.13 BLEU)
- FR→EN: 60.04 vs 54.90 (+5.14 BLEU)

*Note: BASE\_HQ performs better on COMET*





We can build better NMT systems with synthetic data + FM augmentation compared to vanilla NMT.

# Results





## 2. Using reduced monolingual data sets (keeping all bilingual data)



### System Behavior Trends

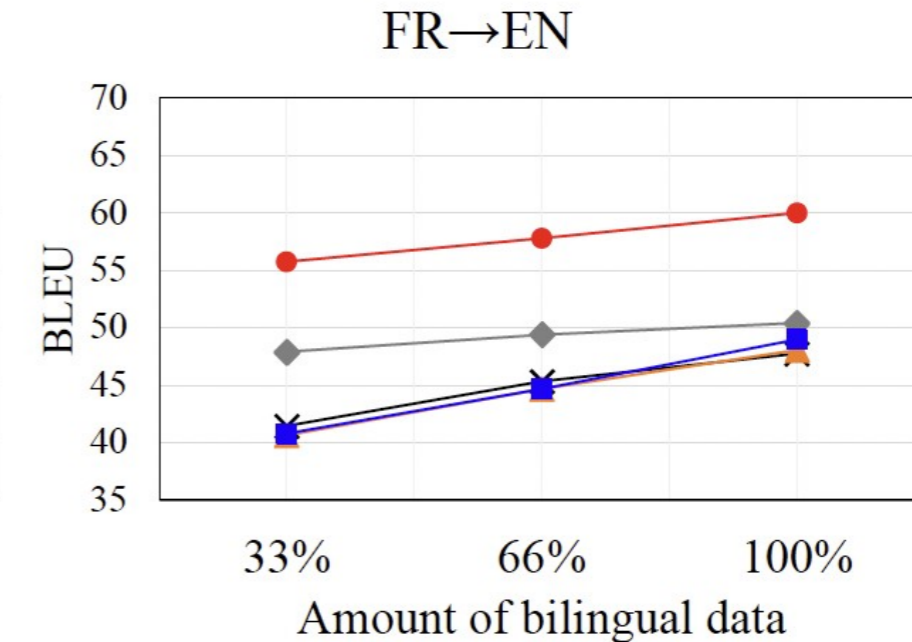
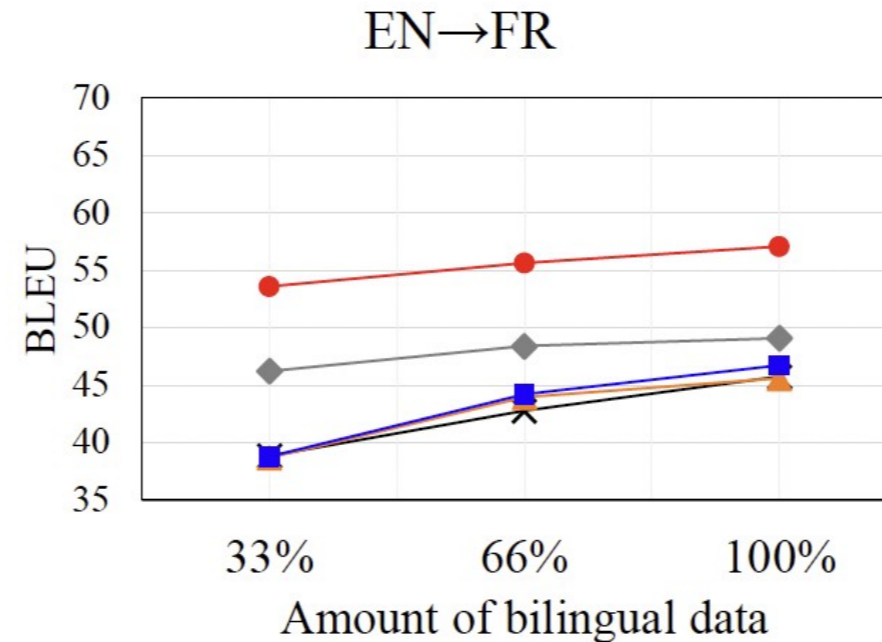
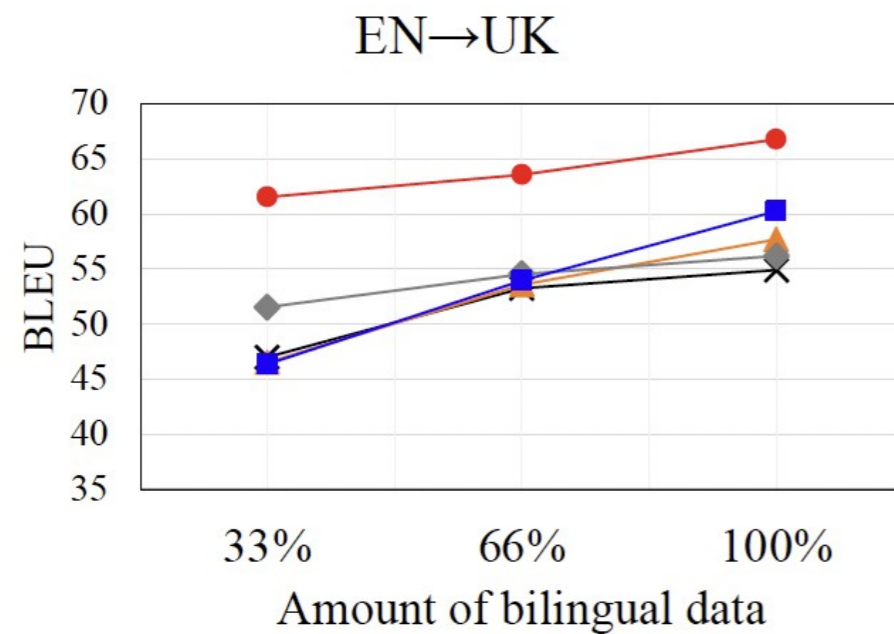
-  **BT** reaches a plateau around 40-60% of data size
-  **NFRmono** keeps slightly improving with additional data
-  Both baselines show diminishing returns with more data
-  Clear performance differences emerge early

### Our System Performance

-  **Outperforms both baselines** in all data conditions
-  **Largest improvements** with maximum available data
-  **Consistent upward trend** across all language pairs
-  **Robust performance** even with limited monolingual data

# Results

## 3. Using reduced bilingual data sets (keeping all monolingual data)



✕ BASE    ▲ NFR    ◆ BT  
■ NFRmono    ● Proposed

✕ BASE    ▲ NFR    ◆ BT  
■ NFRmono    ● Proposed

✕ BASE    ▲ NFR    ◆ BT  
■ NFRmono    ● Proposed



**All baselines improve** with larger bilingual data, confirming expected scaling behavior



**NFRmono benefits most** from larger bilingual data among baseline methods



**Proposed system outperforms** all baselines across all data scenarios



**Large relative improvements** maintained even in lowest resource scenario (33% data)

# Conclusions

- A simple, yet effective approach for improving NMT systems
- Backtranslation + FM-augmentation
  - No changes required to the underlying NMT architecture (data augmentation)
  - Wider adoption than alternatives?
- Additional monolingual data sets in the target language can be more efficient than existing alternatives
- Large (statistically significant) improvements compared to best alternative in all data configurations (including lower resource scenarios)
- Using high-quality, large data sets instead does not improve the results drastically
- Similar trends observed across different metrics

# Limitations & Future Work

## **Limitations:**

- Only 2-domains and 3 language pairs tested
- Only automated evaluation (no human evaluation)

## **Future work:**

- Can similar improvements be observed using FM-augmentation in LLMs?
  - Using LLMs for back-translating monolingual data
  - Integrating FMs through in-context learning
- Can we generate additional monolingual data in the target language, when they are not available?
  - Use LLMs for domain-specific data generation
- Can we even use this approach when only monolingual data sets are available in the target language (no bilingual data)?

# REFERENCES

- Qian Cao and Deyi Xiong. 2018. Encoding gated translation memory into neural machine translation. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, pages 646–657.
- Qiuxiang He, Guoping Huang, Qu Cui, Li Li, and Lemao Liu. 2021. Fast and accurate neural machine translation with translation memory. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 3170–3180, Online. Association for Computational Linguistics.
- Maxime Bouthors, Josep Crego, and François Yvon. 2023. Towards example-based NMT with multi-Levenshtein transformers. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 1830–1846, Singapore. Association for Computational Linguistics.
- Bram Bulté and Arda Tezcan. 2019. Neural fuzzy repair: Integrating fuzzy matches into neural machine translation. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 1800–1809, Florence, Italy. Association for Computational Linguistics.
- Jitao Xu, Josep Crego, and Jean Senellart. 2020. Boosting neural machine translation with similar translations. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 1580–1590, Online. Association for Computational Linguistics.
- Sergey Edunov, Myle Ott, Michael Auli, and David Grangier. 2018. Understanding back-translation at scale. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, pages 489–500, Brussels, Belgium. Association for Computational Linguistics.

# Thomas Moerman

PhD Researcher

LANGUAGE AND TRANSLATION TECHNOLOGY TEAM (LT3)

 Thomas.Moerman@UGent.be

[www.ugent.be](http://www.ugent.be)

<https://lt3.ugent.be/>

Tezcan, A., Skidanova, A., & Moerman, T. (2024). Improving Fuzzy Match Augmented Neural Machine Translation in Specialised Domains through Synthetic Data. *The Prague Bulletin of Mathematical Linguistics*, 122, 9–42.

