

EvolveMT: An Ensemble MT Engine Improving Itself with Usage Only

INTRODUCTION

- Referenceless MT quality estimation metrics like COMET-QE have been shown to outperform human crowd workers in evaluating MT outputs to provide a reliable quality estimation of their hypotheses.
- **Motivation:** Existing MT quality estimation metrics rely on a pre-training scheme with referenced datasets or human evaluations, which limit their adaptability without continuous human-feedback.
- EvolveMT improves translation performance by efficiently combining multiple machine-translation (MT) engines by dynamically selecting the most suitable output for each translation request.
- The method selects the best MT engine for each translation request using online learning and a pre-trained quality estimation metric to supervise itself without the need for reference translations.
- EvolveMT can adapt itself to changes in the domain or the MT engines without requiring retraining and predicts a subset of translation engines to explore based on source sentence features.
- Experimental results show that EvolveMT achieves similar translation accuracy at a lower cost compared to selecting the best translation from all available engines using the MT quality estimator.
- **Contributions:** EvolveMT is the first MT system to adapt itself after deployment for translation requests incoming from the production environment without costly retraining on human feedback.
- **Major Benefits:** Dynamic adaptation to changes without retraining, selective engine calling based on source sentence features, and configurable exploration mechanism of the EvolveMT algorithm.

METHODOLOGY

- EvolveMT is an ensemble MT engine that predicts the best MTs to explore based on the source-only features of each sentence and continuously adapts itself to improve its translation performance.
- EvolveMT uses a teacher-student learning scheme and knowledge distillation and leverages the state-of-the-art neural MT quality metrics and online learning to adapt without human-in-the-loop evaluation.
- The method uses a multi-class classification model to predict the best MT engine for each translation request and utilizes active-learning and contextual bandit learning for the exploration of MTs
 - EvolveMT explores multiple MT engines based on their predicted probabilities and selects the best translation based on COMET-QE scores.
 - Active-learning and contextual bandit learning techniques are employed to prioritize exploration and update the classifier iteratively.
 - The proposed method intelligently ranks translation requests in a queue based on classifier uncertainty and re-ranks after each response.
- The classifier is trained using an online AutoML framework and updated with COMET-QE scores to reinforce the selection process and continuously adapt to the received translation requests.

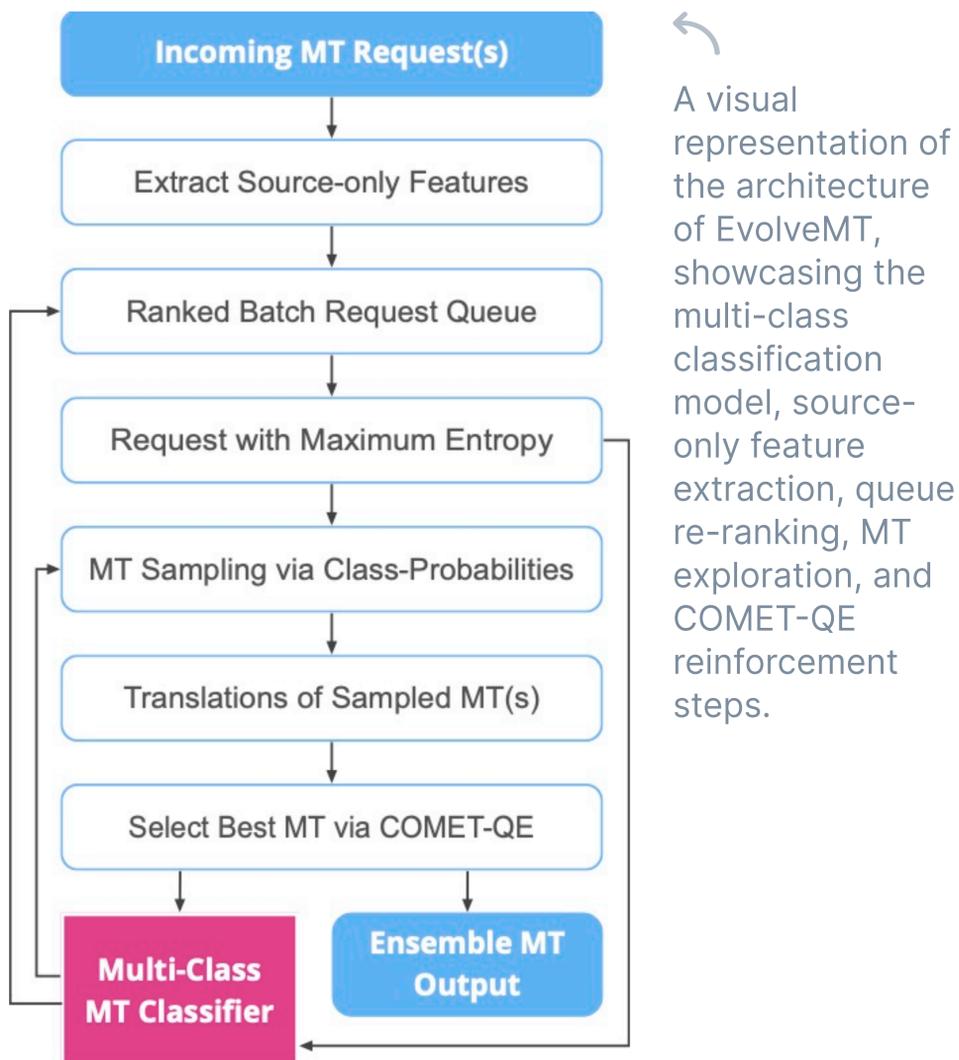
CONCLUSION

- EvolveMT is the first MT system to improve itself after deployment without human feedback and allows continuous adaptation and learning of an MT ensemble without costly retraining.
 - EvolveMT is a novel approach for a multi-system MT (MSMT) engine that improves its translation quality continuously.
 - The method eliminates the need for costly human feedback or retraining, making it more practical and cost-effective.
- EvolveMT outperforms AutoMode and COMET-QE on custom datasets in terms of translation accuracy and cost. The adaptability and efficiency make it a promising approach for multi-system MT.

- Experimental results demonstrate the effectiveness of EvolveMT in achieving accurate translations at a lower cost, leading to significant improvements in translation accuracy and cost-efficiency.
- The proposed method has the potential to enhance MT systems in real-world production environments without requiring extensive manual intervention or model updates.

ARCHITECTURE

The architecture of the proposed method for an ensemble MT engine that improves itself with usage. The multi-class MT classifier, supervised by COMET-QE selections of translations, drives the whole process.



DRIVING ALGORITHM

The pseudo-code of the driving algorithm in EvolveMT, describing the process of selecting the best MT engine, exploration, and continuous learning.

Algorithm 1 EvolveMT with online active-learning

```

Require:  $MTQueue$ : list of tuples (source text, features),
and  $MaxMTs$ : maximum number of MTs to sample
while  $len(MTQueue) > 0$  do
   $Classifier.rankByUncertainty(MTQueue)$ 
   $source, feats \leftarrow MTQueue.popMaxEntropyItem()$ 
   $predMT, classProbs \leftarrow Classifier.predict(feats)$ 
   $predTrans \leftarrow Translate(source, predMT)$ 
   $randMTs \leftarrow sampleMTs(classProbs, MaxMTs)$ 
   $maxEnt \leftarrow normalizedEntropy(classProbs)$ 
  if  $randMTs_0 = predMT$  and  $maxEnt < \alpha$  then
     $Classifier.learn(feats, predMT)$ 
  else
     $sampld \leftarrow Translate(source, randMTs)$ 
     $randScores \leftarrow CometQE(source, sampld)$ 
     $predMTScore \leftarrow CometQE(source, predTrans)$ 
    if  $max(randScores) > predMTScore$  then
       $Classifier.learn(feats, randMTs)$ 
       $IndexOfBestMT \leftarrow randScores.argmax()$ 
       $predTrans \leftarrow sampld_{IndexOfBestMT}$ 
    else
       $Classifier.learn(feats, predMT)$ 
    end if
  end if
   $respondMTRequest((source, predTrans))$ 
end while

```

EXPERIMENTAL RESULTS

The reference-based MT quality scores of the mean and best MT system, AutoMode and COMET-QE ensembles, and EvolveMT ensemble for various MaxMTs parameters (indicated in parentheses next to EvolveMT).

Model	Cost	Δ	English-to-Czech (en-cs)				English-to-German (en-de)				English-to-Russian (en-ru)			
			DA	Δ	HTER	Δ	DA	Δ	HTER	Δ	DA	Δ	HTER	Δ
Mean MT	12.833	-17.20%	0.807	-3.52%	0.139	5.33%	0.550	-2.66%	0.062	3.74%	0.571	4.76%	0.114	-5.11%
Best MT	20.000	29.04%	0.867	3.55%	0.126	-5.12%	0.586	3.71%	0.055	-7.84%	0.617	13.21%	0.104	-13.73%
COMET-QE	77.000	396.81%	0.900	7.50%	0.119	-10.41%	0.605	7.12%	0.053	-11.70%	0.658	20.86%	0.096	-20.27%
EvolveMT (1)	12.312	-20.56%	0.851	1.70%	0.130	-2.08%	0.567	0.43%	0.060	-0.84%	0.605	11.07%	0.108	-10.68%
EvolveMT (2)	23.442	51.25%	0.870	3.99%	0.125	-5.25%	0.586	3.78%	0.057	-5.96%	0.627	15.19%	0.103	-14.69%
EvolveMT (3)	32.358	108.77%	0.878	4.88%	0.123	-6.75%	0.591	4.66%	0.056	-7.62%	0.637	16.95%	0.100	-16.63%
EvolveMT (4)	39.905	157.47%	0.882	5.43%	0.123	-7.30%	0.596	5.40%	0.055	-8.44%	0.643	18.13%	0.099	-17.49%
EvolveMT (5)	46.067	197.23%	0.887	5.94%	0.122	-8.10%	0.598	5.76%	0.055	-9.29%	0.647	18.81%	0.098	-18.20%
EvolveMT (6)	51.095	229.67%	0.887	6.01%	0.122	-8.11%	0.599	6.04%	0.054	-9.54%	0.651	19.43%	0.098	-18.70%
AutoMode	15.499	0.00%	0.837	0.00%	0.132	0.00%	0.565	0.00%	0.060	0.00%	0.545	0.00%	0.120	0.00%

- The MT quality increases as COMET-DA scores increase and COMET-HTER scores decrease. The percentage changes of all scores and costs (per million characters) versus AutoMode ensemble are also indicated.