

# Test-time Adaptation for Machine Translation Evaluation by Uncertainty Minimization



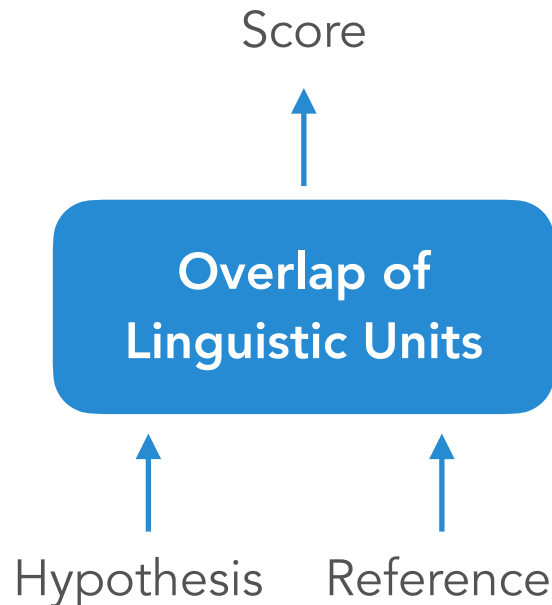
@ACL 2023

Runzhe Zhan, Xuebo Liu, Derek F. Wong, Cuilian Zhang, Lidia S. Chao, Min Zhang

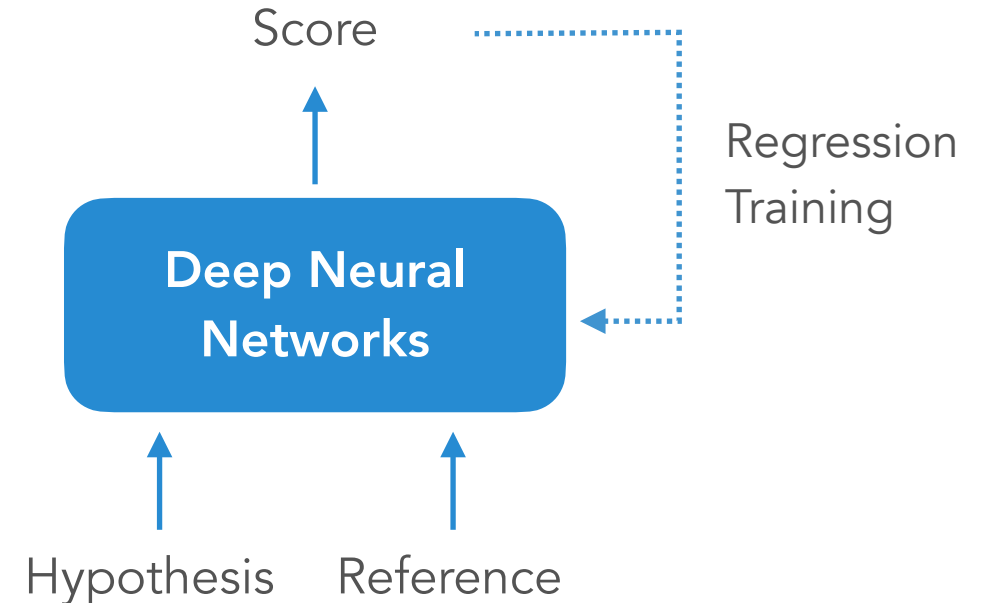
# Introduction

## MT Evaluation

- **Metric:** Automatically quantify the translation quality.
- **Paradigms:** Traditional metrics (e.g., BLEU), Neural metrics (e.g., COMET).



**Traditional Metrics: BLEU, chrF, etc.**

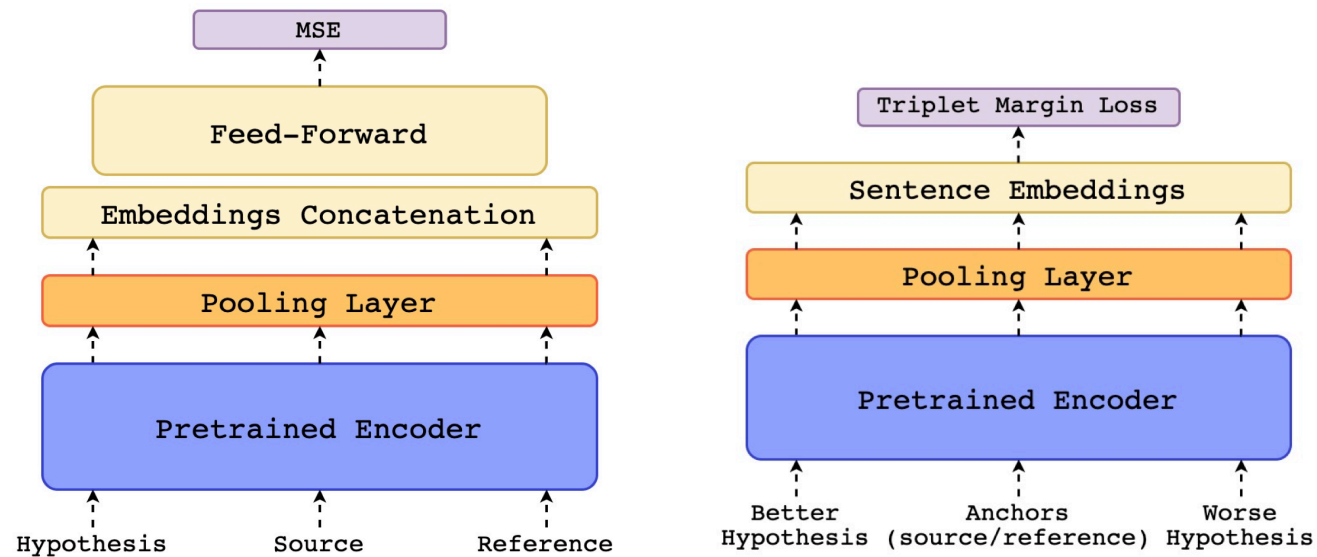


**Neural Metrics: COMET, BLEURT, etc.**

# Introduction

## Related Work

- **Representative Metric - COMET:** Fine-tune XLM-R pre-trained model with human rating data.
- Publicly available rating data solely come from WMT-News domain.

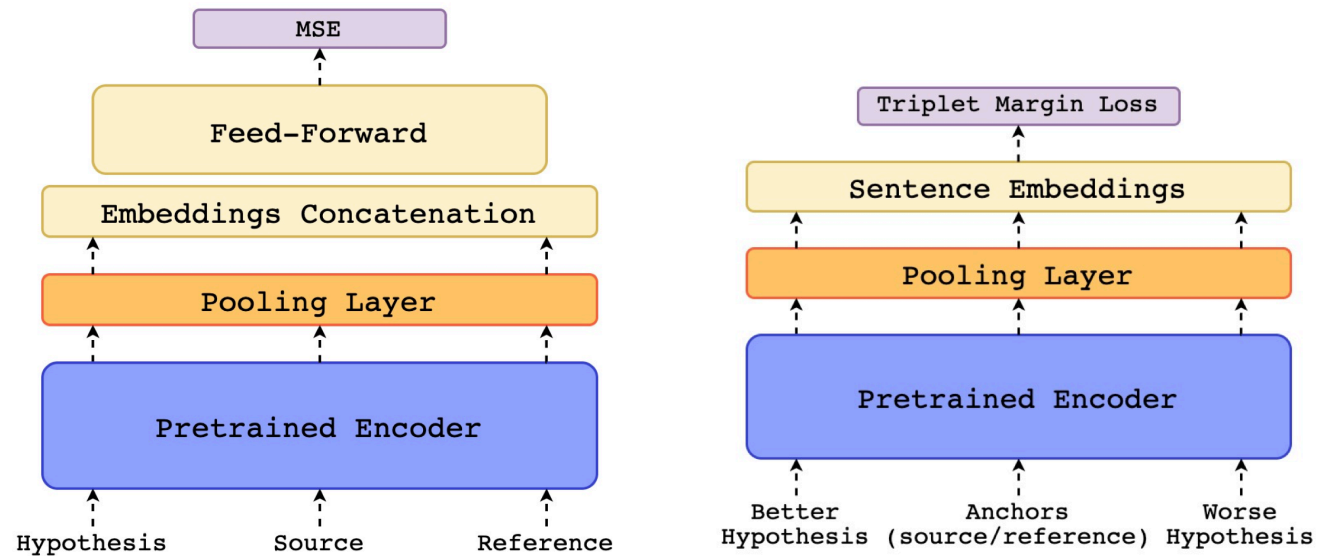


COMET Architectures.

# Introduction

## Related Work

- **Representative Metric - COMET:** Fine-tune XLM-R pre-trained model with human rating data.
- Publicly available rating data solely come from WMT-News domain.
- **Problem:** Neural metrics were trained on WMT-News human rating data.



COMET Architectures.

# Introduction

## Background

- **Problem:** Neural metrics were trained on **WMT-News human rating data**.
- **Potential Risk:** **Robustness problem** when evaluating out-of-distribution (OOD) samples.

# Introduction

## Background

- **Problem:** Neural metrics were trained on **WMT-News human rating data**.
- **Potential Risk:** **Robustness problem** when evaluating out-of-distribution (OOD) samples.
- **Results:** Neural metrics sometimes **underperform** traditional metrics.

	Traditional metrics			Neural metrics		
	BLEU	TER	ChRF	COMET-DA	COMET-QE	COMET-MQM
Chinese-English (TED Domain)	32.4	42.1	36.3	25.1	-20.9	26.6

System-level Pearson correlations (%) of metrics with human scores on WMT21 Metrics Shared Task - MQM.

# Introduction

## Dilemma

- **Direct Solution:** Collect **human scores** for out-of-distribution (OOD) samples.
  - ✗ **Cost:** **Expensive** to collect annotated data!

# Introduction

## Dilemma

- **Direct Solution:** Collect **human scores** for out-of-distribution (OOD) samples.
  - ✗ **Cost:** **Expensive** to collect annotated data!
- **Challenge:** Can we alleviate OOD problem **without collecting annotated data**?
  - **OOD vs. In-domain:** Performance degradation means more prediction errors.



# Introduction

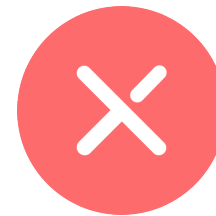
## Dilemma

- **Direct Solution:** Collect **human scores** for out-of-distribution (OOD) samples.
  - ✗ **Cost:** **Expensive** to collect annotated data!
- **Challenge:** Can we alleviate OOD problem **without collecting annotated data**?
  - **OOD vs. In-domain:** Performance degradation means more prediction errors.
    - 🤔 **What factors are related to the prediction errors from model perspective?**

# Introduction

## Dilemma

- **Direct Solution:** Collect **human scores** for out-of-distribution (OOD) samples.
    - ✗ **Cost:** **Expensive** to collect annotated data!
  - **Challenge:** Can we alleviate OOD problem **without collecting annotated data**?
    - **OOD vs. In-domain:** Performance degradation means more prediction errors.
- 🤔 **What factors are related to the prediction errors from model perspective?**



Prediction Error

💡 **Minimize [?] → Minimize (model's) prediction error**

Observation

# Introduction

# Introduction

## Observation

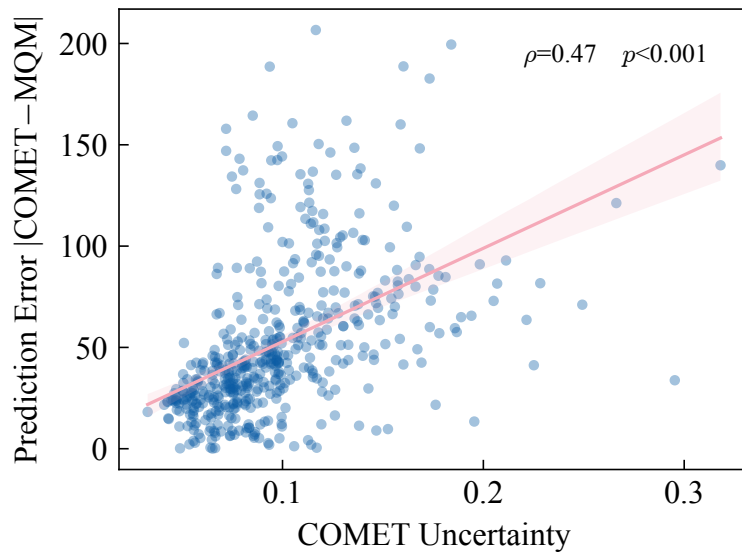
🤔 What factors are related to the prediction errors from model perspective?

# Introduction

## Observation

🤔 What factors are related to the prediction errors from model perspective?

- **Model uncertainty** reflects the risk of model's prediction.
- **Observation:** Model uncertainty **positively aligns** with its prediction errors.
  - Also observed by Glushkova et al. (2021).



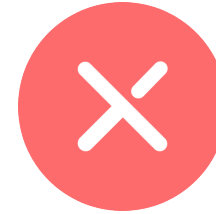
MQM: Human Ratings  
 COMET: A Neural Metric  
 Prediction Error: Difference between metric score and human score.

# Introduction

## Motivation



Model Uncertainty

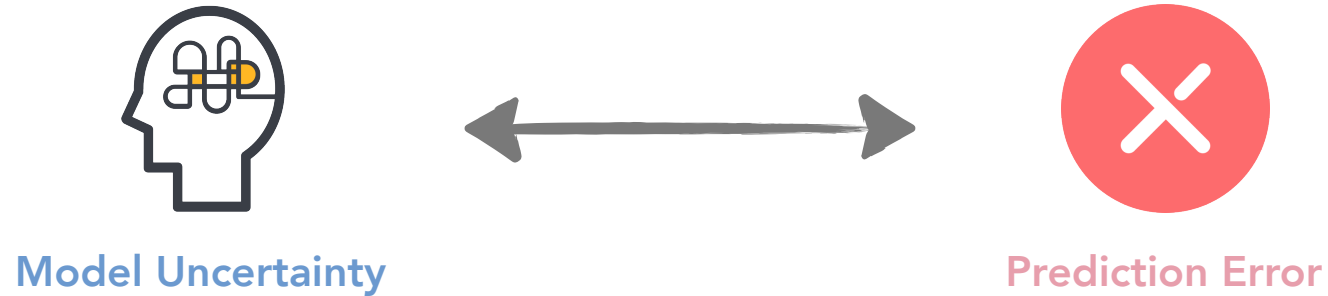


Prediction Error

🤔 **Key Idea: Minimize (model's) uncertainty → Minimize (model's) prediction error**

# Introduction

## Motivation

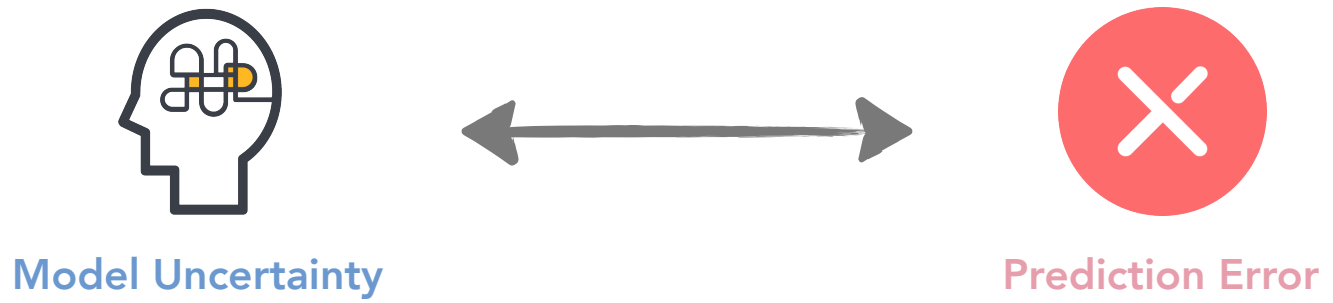


🤔 **Key Idea: Minimize (model's) uncertainty → Minimize (model's) prediction error**

- ✓ Can be estimated **during test time**
- ✓ Can be estimated **without additional data**
- ✓ Make the model correct the predictions by itself

# Introduction

## Motivation



🤔 **Key Idea: Minimize (model's) uncertainty → Minimize (model's) prediction error**

- ✓ Can be estimated **during test time**
- ✓ Can be estimated **without additional data**
- ✓ Make the model correct the predictions by itself



**Test-time Adaptation by  
Uncertainty Minimization  
(TaU)**



# Methodology

## Proposal: TaU

### 🤔 Key Research Questions:

- 1) How can we **estimate** the uncertainty **for metrics' model** ?
- 2) How can we **reduce** the uncertainty by **test-time adaptation** ?

# Methodology

## Proposal: TaU

### 🤔 Key Research Questions:

- 1) How can we **estimate** the uncertainty **for metrics' model** ?
- 2) How can we **reduce** the uncertainty by **test-time adaptation** ?

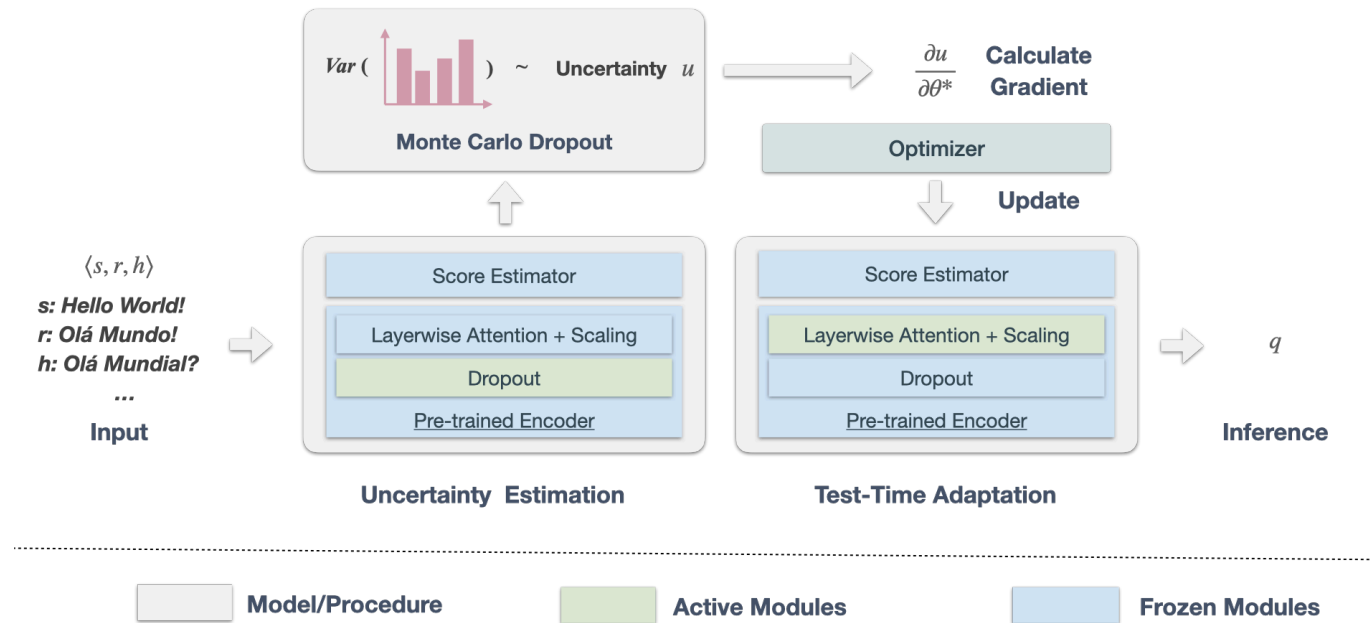


Illustration of the proposed method: Test-time Adaption by Uncertainty estimation (TaU)

# Methodology

## Proposal: TaU

### 🤔 Key Research Questions:

- 1) How can we **estimate** the uncertainty **for metrics' model** ?
- 2) How can we **reduce** the uncertainty by **test-time adaptation** ?

### Uncertainty Estimation (RQ.1)

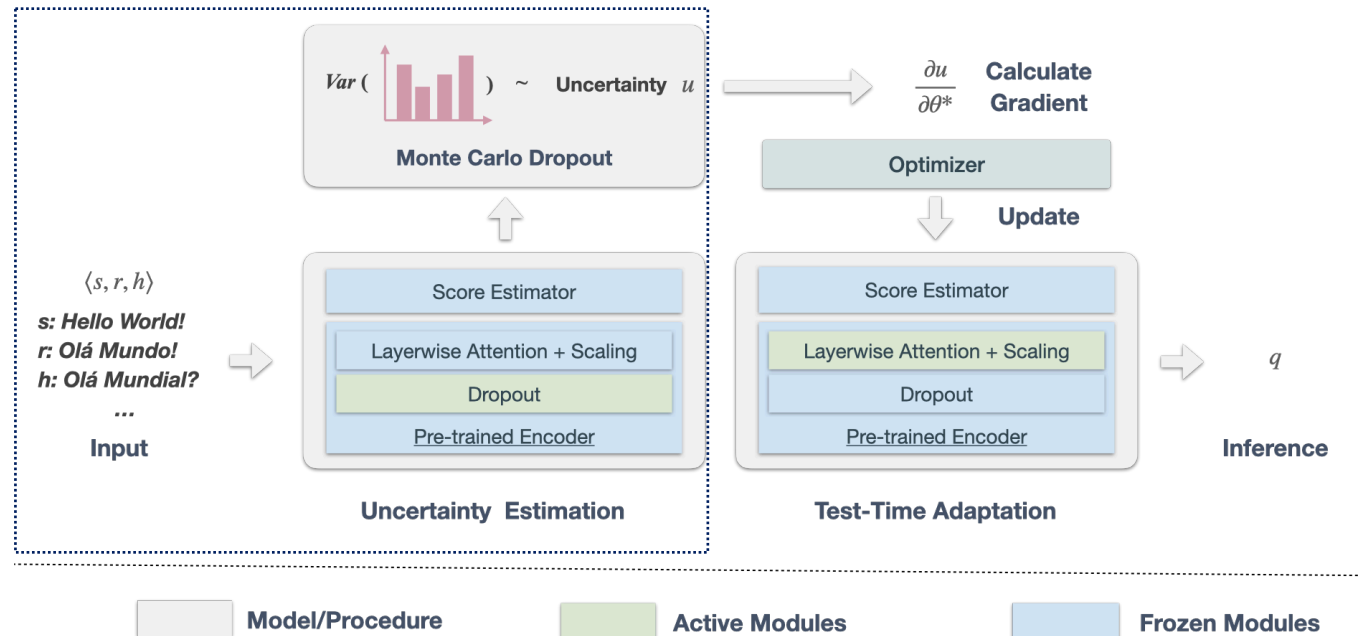


Illustration of the proposed method: Test-time Adaption by Uncertainty estimation (TaU)

# Methodology

## Proposal: TaU

### 🤔 Key Research Questions:

- 1) How can we **estimate** the uncertainty **for metrics' model** ?
- 2) How can we **reduce** the uncertainty by **test-time adaptation** ?

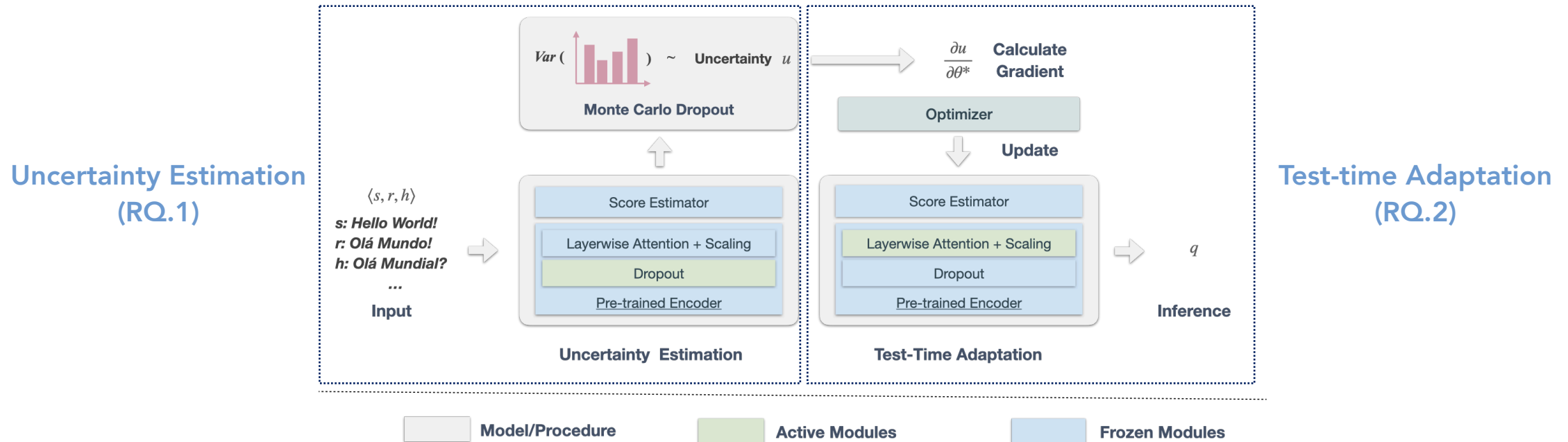
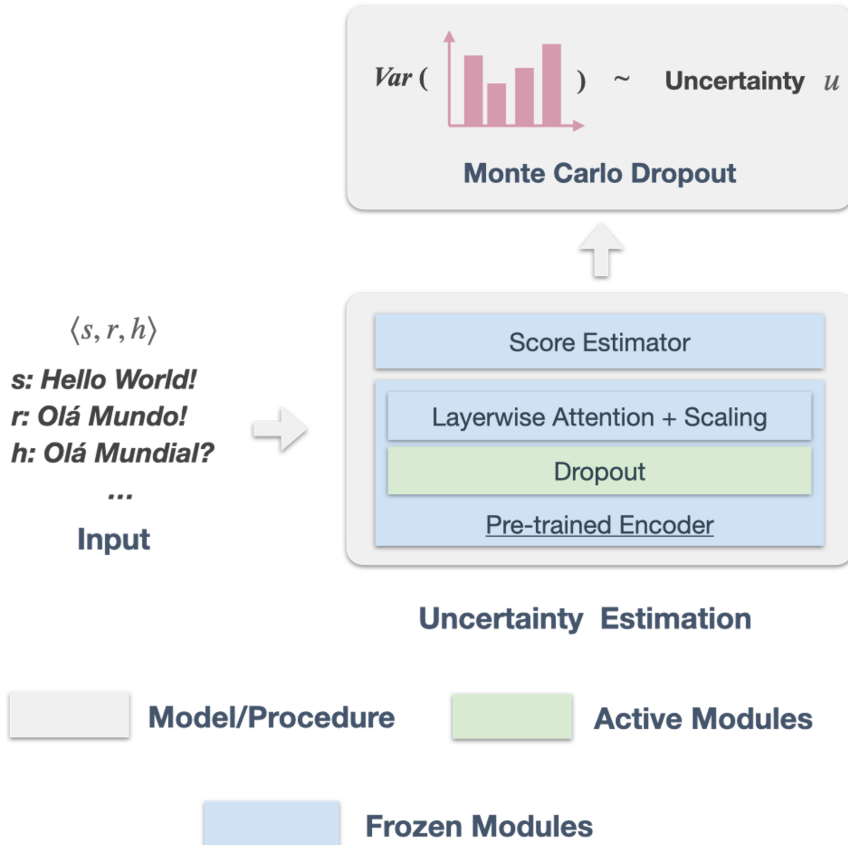


Illustration of the proposed method: Test-time Adaption by Uncertainty estimation (TaU)

# Methodology

## Proposal: TaU

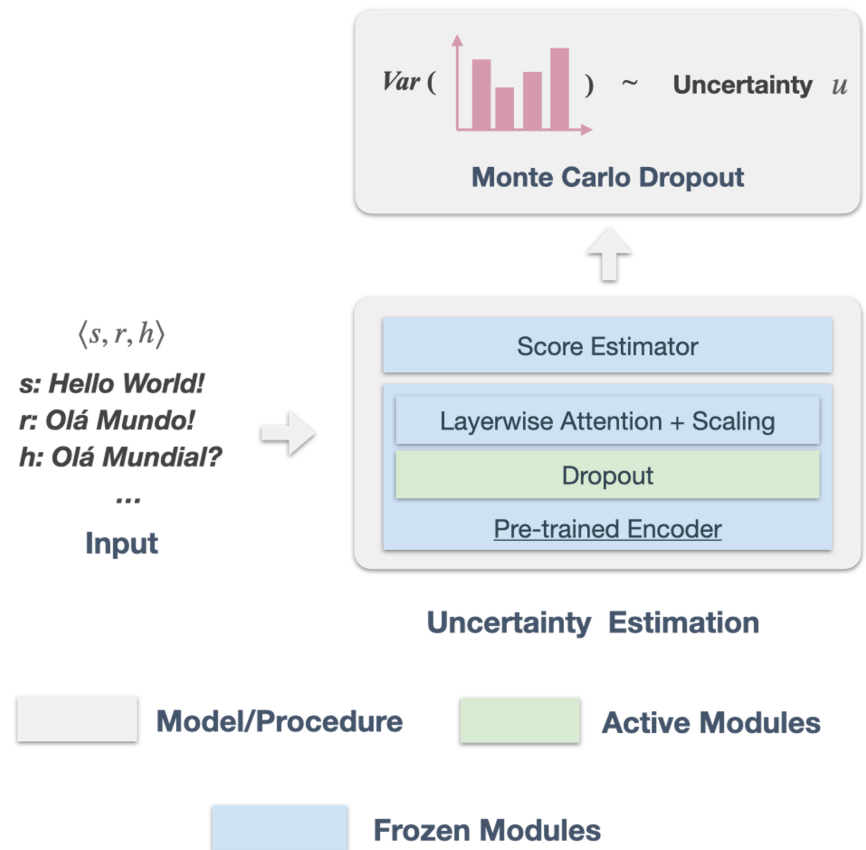


- **Uncertainty Estimation**

- Uncertainty = Variance of **scoring distribution**
- Regression model only provides **a single score**

# Methodology

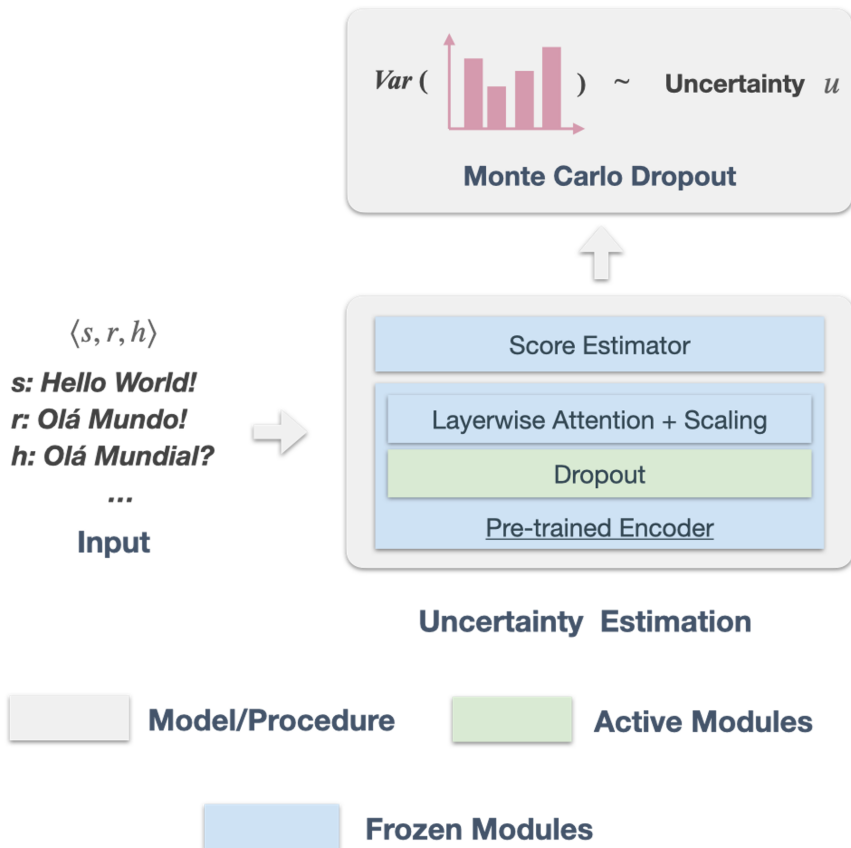
## Proposal: TaU



- **Uncertainty Estimation**
  - Uncertainty = Variance of **scoring distribution**
  - Regression model only provides **a single score**
- **Monte Carlo** sampling method:
  - Randomly active some **dropout layers** and perform K-times forward-propagation.

# Methodology

## Proposal: TaU



- **Uncertainty Estimation**

- Uncertainty = Variance of **scoring distribution**
- Regression model only provides **a single score**

- **Monte Carlo** sampling method:

- Randomly active some **dropout layers** and perform K-times forward-propagation.

- Uncertainty = **Variance of K-times prediction**

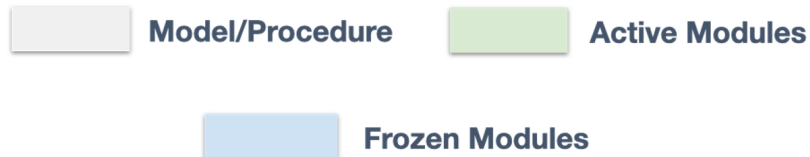
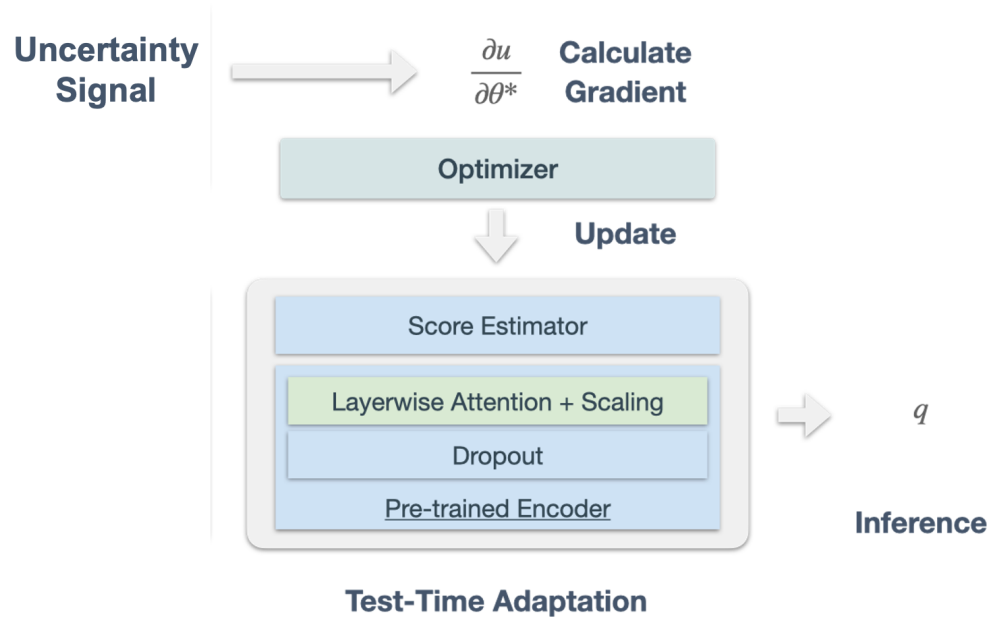
$$u(\langle h, s, \cdot \rangle) = \mathbf{Var}(\{M(\langle h, s, \cdot \rangle; \theta_k)\}_{k=1}^K)$$

↑  
Input Data

↑  
Metric Model (w/ Dropout)

# Methodology

## Proposal: TaU



- **Test-time adaptation**

- Online optimization
- Objective function: **minimize the uncertainty**

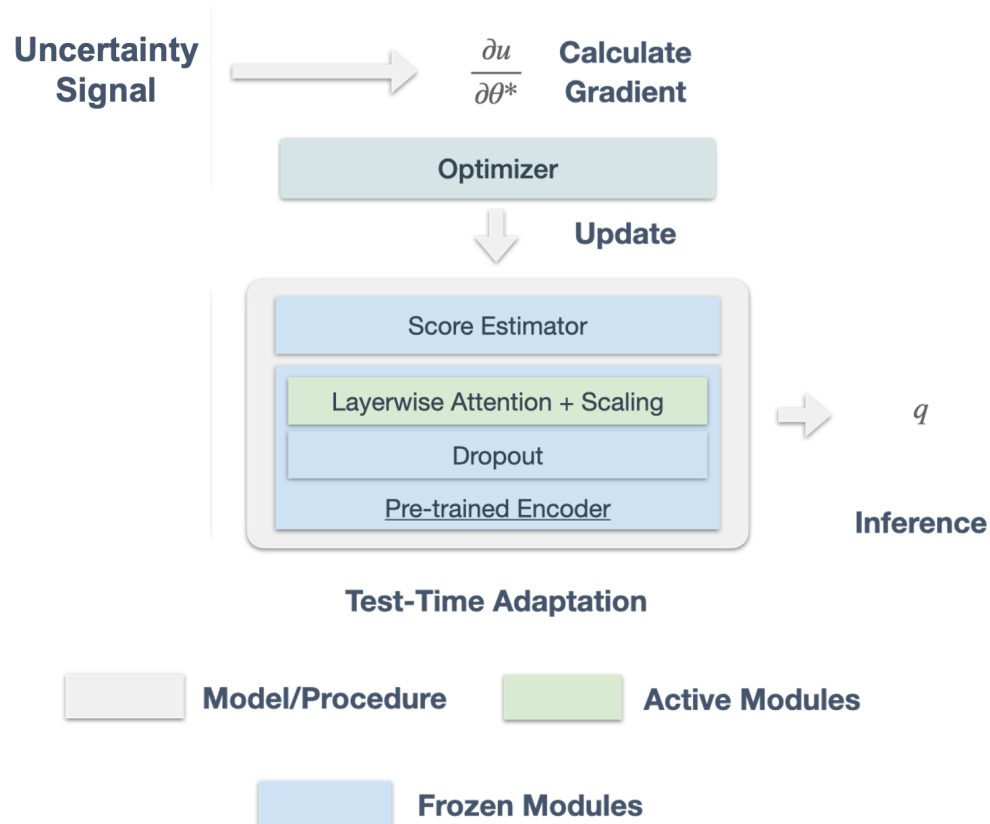
$$\theta^* = \arg \min_{\theta^*} \mathbb{E}_{\langle h, s, \cdot \rangle \in \mathcal{D}} [u(\langle h, s, \cdot \rangle)]$$

↑  
Optimization of partial modules



# Methodology

## Proposal: TaU



- **Test-time adaptation**

- Online optimization
- Objective function: **minimize the uncertainty**

$$\theta^* = \arg \min_{\theta^*} \mathbb{E}_{\langle h, s, \cdot \rangle \in \mathcal{D}} [u(\langle h, s, \cdot \rangle)]$$

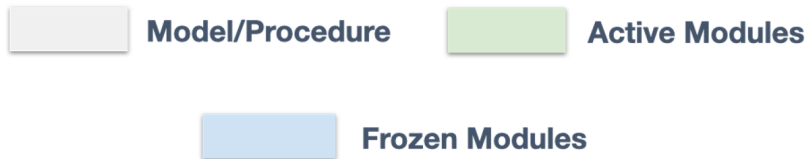
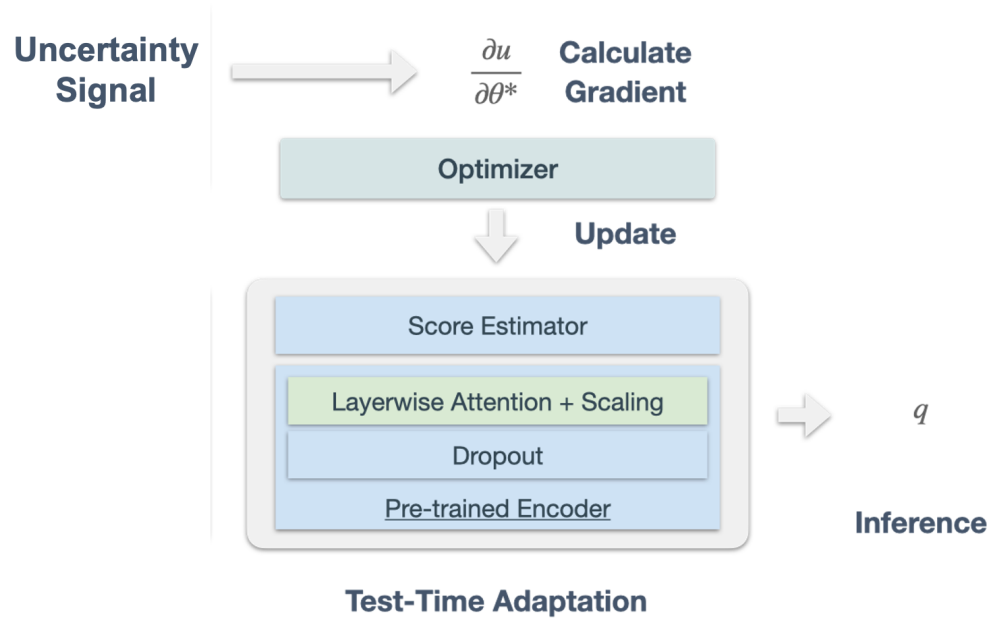
↑  
Optimization of partial modules

- **Choice of Optimization Parameters**

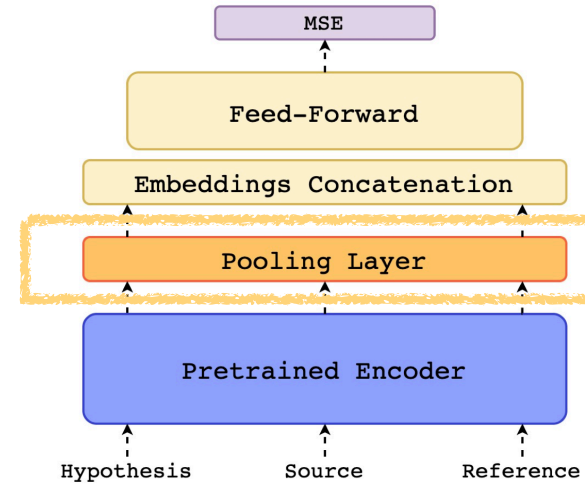
- Do not deviate far from original parameters.
- Only optimize partial parameters.

# Methodology

## Proposal: TaU



- Choice of Optimization Parameters

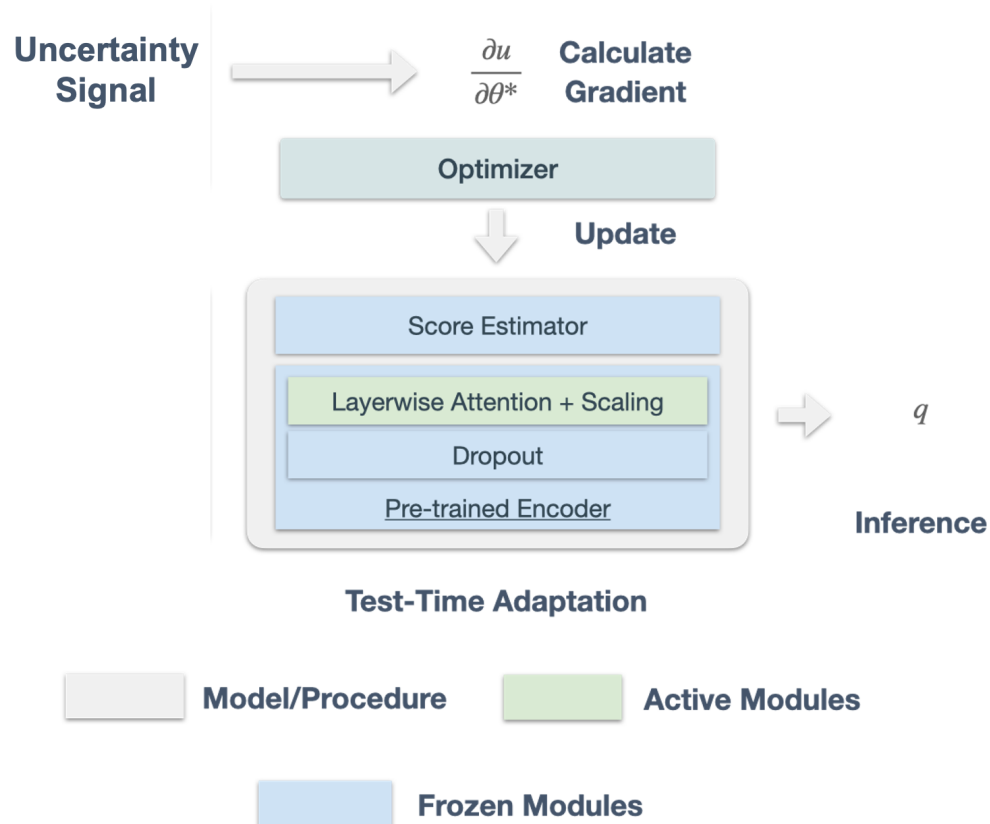


$$\mathbf{O}_{\text{embed}} = \gamma \cdot \sum_{l=1}^L w_l \cdot \text{LayerNorm}(\mathbf{h}_l)$$

↑  
Output of each layer

# Methodology

## Proposal: TaU



- Choice of Optimization Parameters

Domain	LAtt.	LN.	Estim.	$\rho$	Acc.
News	✓	✗	✗	<b>85.7</b>	<b>89.7</b>
	✗	✓	✗	79.5	76.9
	✗	✗	✓	78.7	80.8
	✓	✓	✗	79.6	76.9
	✓	✗	✓	78.6	80.8
TED	✓	✓	✓	78.0	79.4
	✓	✗	✗	<b>85.9</b>	<b>85.9</b>
	✗	✓	✗	79.4	82.1
	✗	✗	✓	77.2	76.9
	✓	✓	✗	79.4	82.1
	✓	✗	✓	77.1	76.9
	✓	✓	✓	76.9	76.9

Ablation Results.

LAtt. = Layerwise Attention | LN. = Layer Normalization

Estim. = Score Estimator

# Methodology

## Algorithm: TaU

**Require:** Model  $\theta$ , System-level evaluation tuple  $\mathcal{D} = \{\langle \mathbf{h}, \mathbf{s}, \cdot \rangle\}$ , Adaptation rate  $\alpha$ , Adaptation times  $J$ .

- 1: Backup original model  $\theta' \leftarrow \theta$
- 2: Select parameters for adaptation  $|\theta^*| \ll |\theta|$
- 3: **for** adaptation iteration  $j = 1, \dots, J$  **do**
- 4:   Score set  $\mathbf{q} = \{\emptyset\}$
- 5:   **for** mini-batch  $\{\langle h, s, \cdot \rangle\}_{i=1}^N \in \mathcal{D}$  **do**
- 6:     Estimate uncertainty  $u$  by Equation 3
- 7:     Optimize  $\theta^* \leftarrow \theta^* - \alpha \nabla_{\theta^*} \frac{1}{N} \sum_{i=1}^N u_i$
- 8:   **end for**
- 9:   Infer score  $[q]$  by Equation 7
- 10:    $\mathbf{q} \leftarrow [q]$
- 11: **end for**
- 12: Restore to original model  $\theta \leftarrow \theta'$
- 13: **return**  $\mathbf{q}$

- **Three steps: Estimate, Adapt, Predict**

### 1. Estimate the model uncertainty

$$u(\langle h, s, \cdot \rangle) = \mathbf{Var}(\{M(\langle h, s, \cdot \rangle; \theta_k)\}_{k=1}^K)$$

$$\mathbf{Var}(P) = \sqrt{\mathbb{E}[(P - \mu_P)^2]}$$

# Methodology

## Algorithm: TaU

**Require:** Model  $\theta$ , System-level evaluation tuple  $\mathcal{D} = \{\langle h, s, \cdot \rangle\}$ , Adaptation rate  $\alpha$ , Adaptation times  $J$ .

- 1: Backup original model  $\theta' \leftarrow \theta$
- 2: Select parameters for adaptation  $|\theta^*| \ll |\theta|$
- 3: **for** adaptation iteration  $j = 1, \dots, J$  **do**
- 4:   Score set  $\mathbf{q} = \{\emptyset\}$
- 5:   **for** mini-batch  $\{\langle h, s, \cdot \rangle\}_{i=1}^N \in \mathcal{D}$  **do**
- 6:     Estimate uncertainty  $u$  by Equation 3
- 7:     Optimize  $\theta^* \leftarrow \theta^* - \alpha \nabla_{\theta^*} \frac{1}{N} \sum_{i=1}^N u_i$
- 8:   **end for**
- 9:   Infer score  $[q]$  by Equation 7
- 10:    $\mathbf{q} \leftarrow [q]$
- 11: **end for**
- 12: Restore to original model  $\theta \leftarrow \theta'$
- 13: **return**  $\mathbf{q}$

- **Three steps: Estimate, Adapt, Predict**

### 1. Estimate the model uncertainty

$$u(\langle h, s, \cdot \rangle) = \mathbf{Var}(\{M(\langle h, s, \cdot \rangle; \theta_k)\}_{k=1}^K)$$

$$\mathbf{Var}(P) = \sqrt{\mathbb{E}[(P - \mu_P)^2]}$$

### 2. Adapt by minimizing the uncertainty

$$\theta^* = \arg \min_{\theta^*} \mathbb{E}_{\langle h, s, \cdot \rangle \in \mathcal{D}} [u(\langle h, s, \cdot \rangle)]$$

# Methodology

## Algorithm: TaU

**Require:** Model  $\theta$ , System-level evaluation tuple  $\mathcal{D} = \{\langle h, s, \cdot \rangle\}$ , Adaptation rate  $\alpha$ , Adaptation times  $J$ .

- 1: Backup original model  $\theta' \leftarrow \theta$
- 2: Select parameters for adaptation  $|\theta^*| \ll |\theta|$
- 3: **for** adaptation iteration  $j = 1, \dots, J$  **do**
- 4:   Score set  $\mathbf{q} = \{\emptyset\}$
- 5:   **for** mini-batch  $\{\langle h, s, \cdot \rangle\}_{i=1}^N \in \mathcal{D}$  **do**
- 6:     Estimate uncertainty  $u$  by Equation 3
- 7:     Optimize  $\theta^* \leftarrow \theta^* - \alpha \nabla_{\theta^*} \frac{1}{N} \sum_{i=1}^N u_i$
- 8:   **end for**
- 9:   Infer score  $[q]$  by Equation 7
- 10:    $\mathbf{q} \leftarrow [q]$
- 11: **end for**
- 12: Restore to original model  $\theta \leftarrow \theta'$
- 13: **return**  $\mathbf{q}$

- **Three steps: Estimate, Adapt, Predict**

### 1. Estimate the model uncertainty

$$u(\langle h, s, \cdot \rangle) = \mathbf{Var}(\{M(\langle h, s, \cdot \rangle; \theta_k)\}_{k=1}^K)$$

$$\mathbf{Var}(P) = \sqrt{\mathbb{E}[(P - \mu_P)^2]}$$

### 2. Adapt by minimizing the uncertainty

$$\theta^* = \arg \min_{\theta^*} \mathbb{E}_{\langle h, s, \cdot \rangle \in \mathcal{D}} [u(\langle h, s, \cdot \rangle)]$$

### 3. Predict with the adapted parameters

$$q = M_{\theta + \Delta \theta^*}(\{\langle h, s, \cdot \rangle\})$$

# Experiments

## TaU: Performance

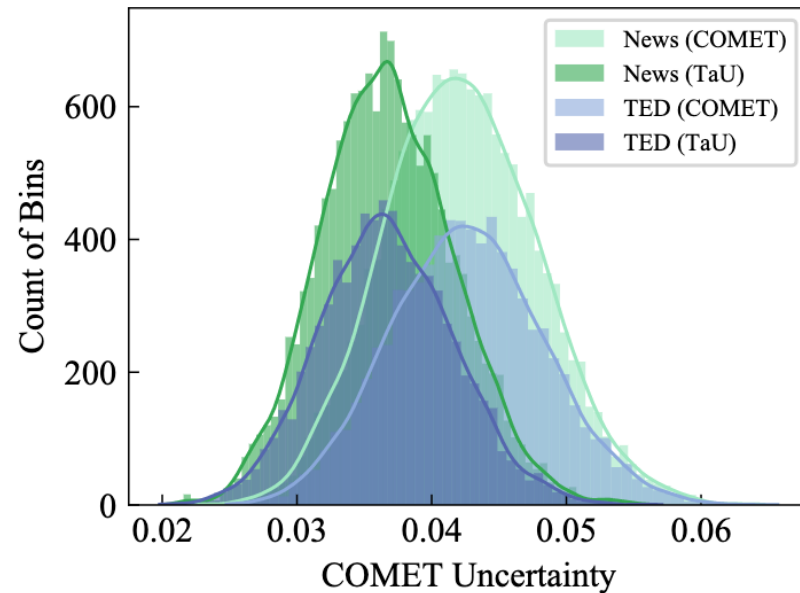
- **Testbed:** COMET models, Developmental data: WMT20
- Improved **system-level Pearson's correlation** on partial multi-domain evaluation tasks.

Metrics	News w/o HT			News w/ HT			TED			Avg.
	En-De	Zh-En	En-Ru	En-De	Zh-En	En-Ru	En-De	Zh-En	En-Ru	
<i>Baselines</i>										
TER	93.0	41.6	-4.1	7.4	-8.5	-28.9	50.6	42.1	69.7	29.2
BLEU	93.7	31.0	50.7	13.2	-15.2	-4.3	62.0	32.4	82.8	38.5
CHRf	89.8	30.2	78.3	1.7	-14.3	12.3	47.1	36.3	82.5	40.4
BERTSCORE	93.0	54.2	62.9	7.4	9.5	-12.3	50.6	30.6	83.1	42.1
COMET-DA <sub>2020</sub>	81.4	51.1	67.6	65.8	22.1	55.6	78.8	25.1	85.9	59.3
COMET-MQM-QE <sub>2021</sub>	71.1	52.9	63.2	79.2	61.9	68.1	69.4	-20.9	88.4	59.3
COMET-MQM <sub>2021</sub>	77.1	62.8	65.9	72.0	33.6	68.5	81.8	26.6	84.1	63.6
<i>Reproduced Results and Our Methods</i>										
◇ COMET-DA <sub>2020</sub>	81.5	51.1	67.5	58.0	26.4	56.8	78.8	25.0	85.9	59.0
+TAU	<b>85.7</b>	<b>53.5</b>	<b>71.0</b>	48.0	<b>27.4</b>	54.5	<b>85.9</b>	<b>28.3</b>	<b>87.3</b>	<b>60.2</b>
◇ COMET-MQM-QE <sub>2021</sub>	71.2	53.0	68.8	79.2	61.9	68.1	69.4	-20.8	81.7	59.2
+TAU	62.8	<b>57.4</b>	<b>70.3</b>	72.0	<b>65.2</b>	<b>78.1</b>	<b>82.9</b>	<b>25.7</b>	80.7	<b>66.1</b>
◇ COMET-MQM <sub>2021</sub>	77.2	62.8	65.9	69.8	48.7	69.7	81.8	26.6	84.1	65.2
+TAU	76.5	<b>69.2</b>	<b>67.2</b>	<b>75.4</b>	<b>67.8</b>	<b>71.5</b>	<b>87.5</b>	24.5	<b>84.9</b>	<b>69.4</b>

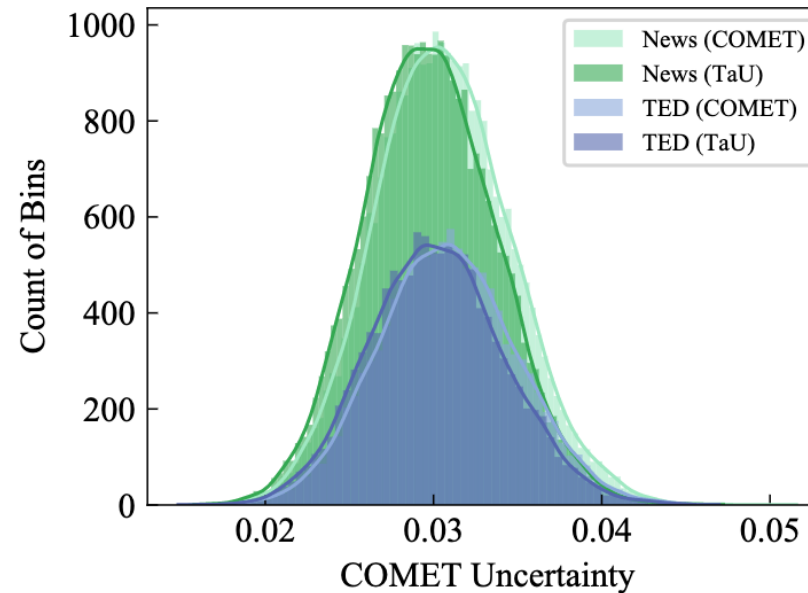
# Analysis

## TaU: Effectiveness

- **Research Goal:** Reduce the uncertainty of OOD samples
- **Validity:** improve the correlation, also reduce the uncertainty



(a) English-German (COMET-DA<sub>2020</sub>)



(b) Chinese-English (COMET-MQM<sub>2021</sub>)

Uncertainty distribution of COMET baselines and corresponding models optimized by TaU



## Conclusions

# Conclusions

# Conclusions

## Conclusions

- **OOD samples** may hinder the application scope of **neural metrics**.

# Conclusions

## Conclusions

- **OOD samples** may hinder the application scope of **neural metrics**.
- We confirmed and also observed the **uncertainty-error relationship** for metric models.
- We propose a **test-time adaptation** method to **reduce inference uncertainty**.

# Conclusions

## Conclusions

- **OOD samples** may hinder the application scope of **neural metrics**.
- We confirmed and also observed the **uncertainty-error relationship** for metric models.
- We propose a **test-time adaptation** method to **reduce inference uncertainty**.
- **OOD evaluation / adaptation** are potential topics for the **Large Language Models**.

# Conclusions

## Conclusions

- **OOD samples** may hinder the application scope of **neural metrics**.
- We confirmed and also observed the **uncertainty-error relationship** for metric models.
- We propose a **test-time adaptation** method to **reduce inference uncertainty**.
- **OOD evaluation / adaptation** are potential topics for the **Large Language Models**.

## Limitations

- **Segment-level** correlation performance is not satisfactory.
- Hyper-parameter searching is still **time-consuming**.
- Cannot fix the errors related to **unseen knowledge**.

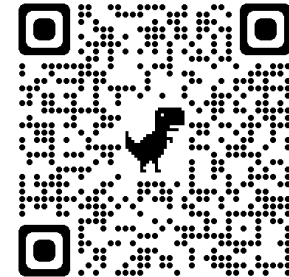
# Conclusions

## Conclusions

- **OOD samples** may hinder the application scope of **neural metrics**.
- We confirmed and also observed the **uncertainty-error relationship** for metric models.
- We propose a **test-time adaptation** method to **reduce inference uncertainty**.
- **OOD evaluation / adaptation** are potential topics for the **Large Language Models**.

## Limitations

- **Segment-level** correlation performance is not satisfactory.
- Hyper-parameter searching is still **time-consuming**.
- Cannot fix the errors related to **unseen knowledge**.



Poster

Thank you for listening!

# Supplementary

## Low-Resource Languages (In-domain)

- **Multi-domain benchmark for MT evaluation is scarce.**
- Experimental results on previous WMT-News benchmark with the same learning rate (did not tune on developmental data).

	Pl-En	Ru-En	Ta-En	Zh-En	En-Pl	En-Ru	En-Ta	En-Zh
COMET-DA	34.5	83.6	0.764	93.1	80.0	92.5	79.8	0.7
+TaU	<b>34.6</b>	<b>84.0</b>	<b>0.774</b>	<b>93.4</b>	79.0	91.6	75.3	<b>1.2</b>

System-level Pearson correlations (%) of metrics with human scores on WMT20 Metrics Shared Task (News).

# Supplementary

## Low-Resource Languages (In-domain)

- Multi-domain benchmark for MT evaluation is scarce.
- Experimental results on previous WMT-News benchmark with the same learning rate (did not tune on developmental data).

	into English				from English			
	Pl-En	Ru-En	Ta-En	Zh-En	En-Pl	En-Ru	En-Ta	En-Zh
COMET-DA	34.5	83.6	0.764	93.1	80.0	92.5	79.8	0.7
+TaU	<b>34.6</b>	<b>84.0</b>	<b>0.774</b>	<b>93.4</b>	79.0	91.6	75.3	<b>1.2</b>

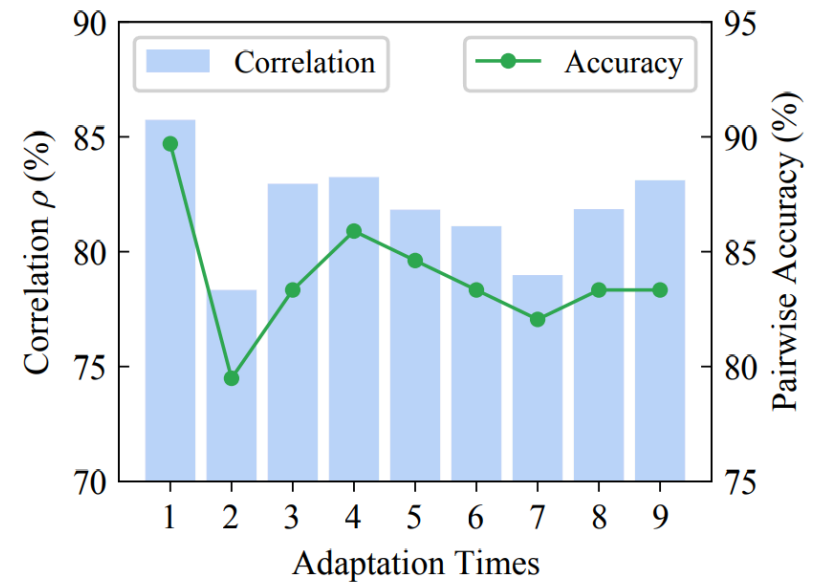
System-level Pearson correlations (%) of metrics with human scores on WMT20 Metrics Shared Task (News).



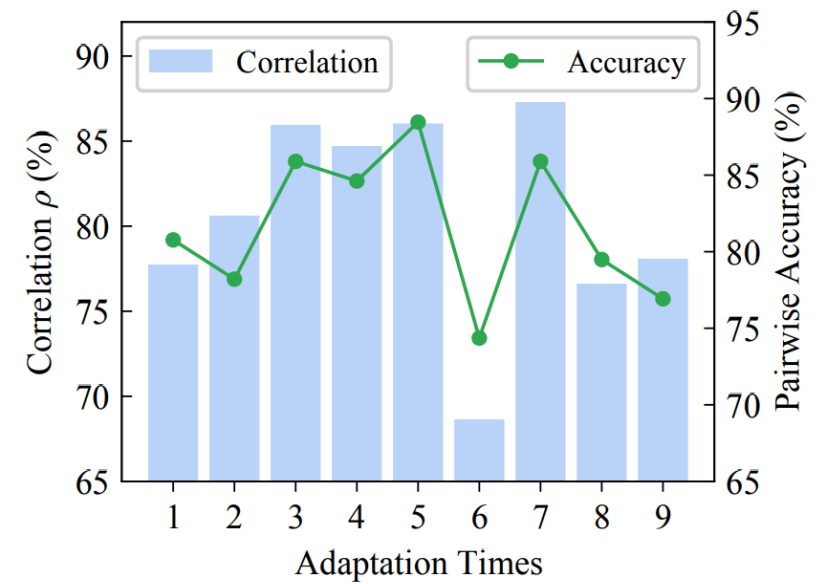
# Supplementary

## Multi-turn Adaptation

- The out-of-distribution data requires more adaptation times than in-domain data, and both of them would suffer from extreme settings.



(a) News



(b) TED

Performance of TAU with different settings of adaptation times.