

Emmanouil Antonios Platanios
e.a.platanios@cs.cmu.edu

Otilia Stretcu
ostretcu@cs.cmu.edu

Graham Neubig
gneubig@cs.cmu.edu

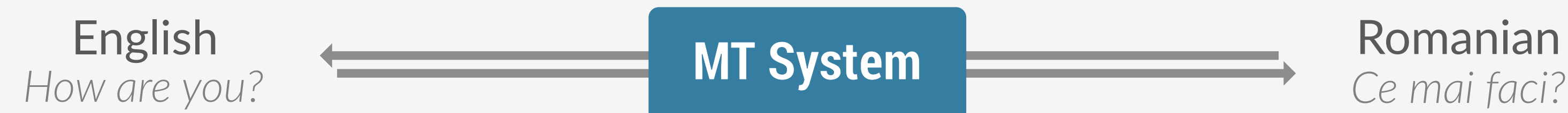
Barnabas Poczos
bapoczos@cs.cmu.edu

Tom M. Mitchell
tom.mitchell@cs.cmu.edu

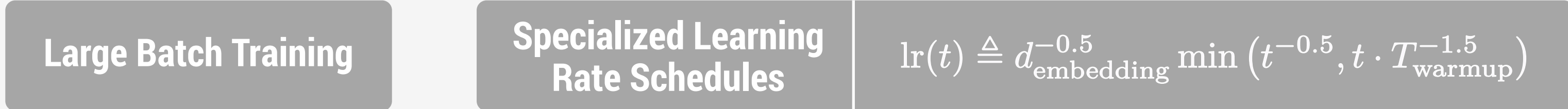
Motivation

MACHINE TRANSLATION

Translate from one language to another:



Large scale neural MT systems are hard to train. For example, Transformers require:



CURRICULUM LEARNING



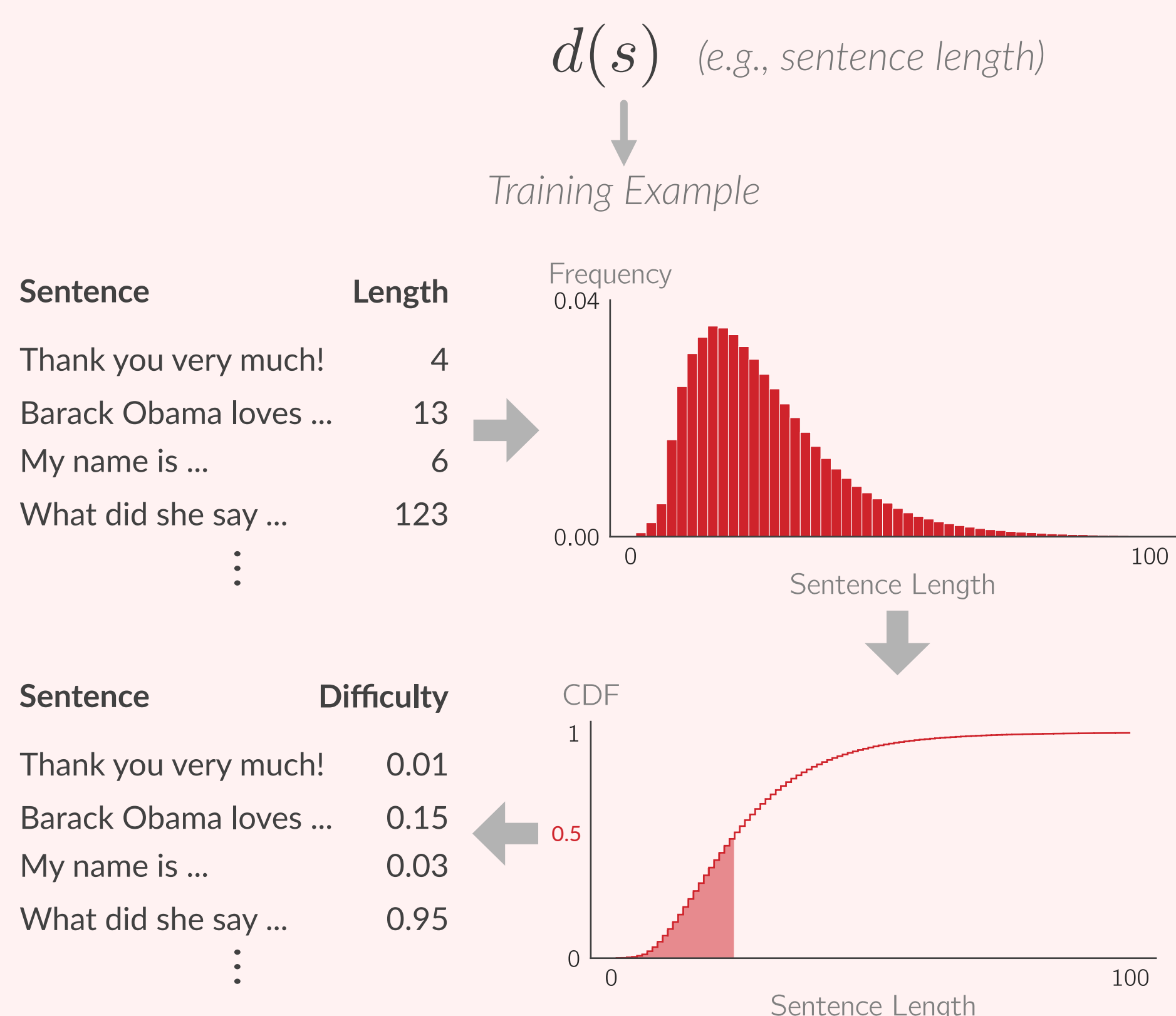
Previous curriculum learning approaches for NMT:



Proposed Approach

DIFFICULTY

Represents the difficulty of a training example that may depend on the state of the learner:

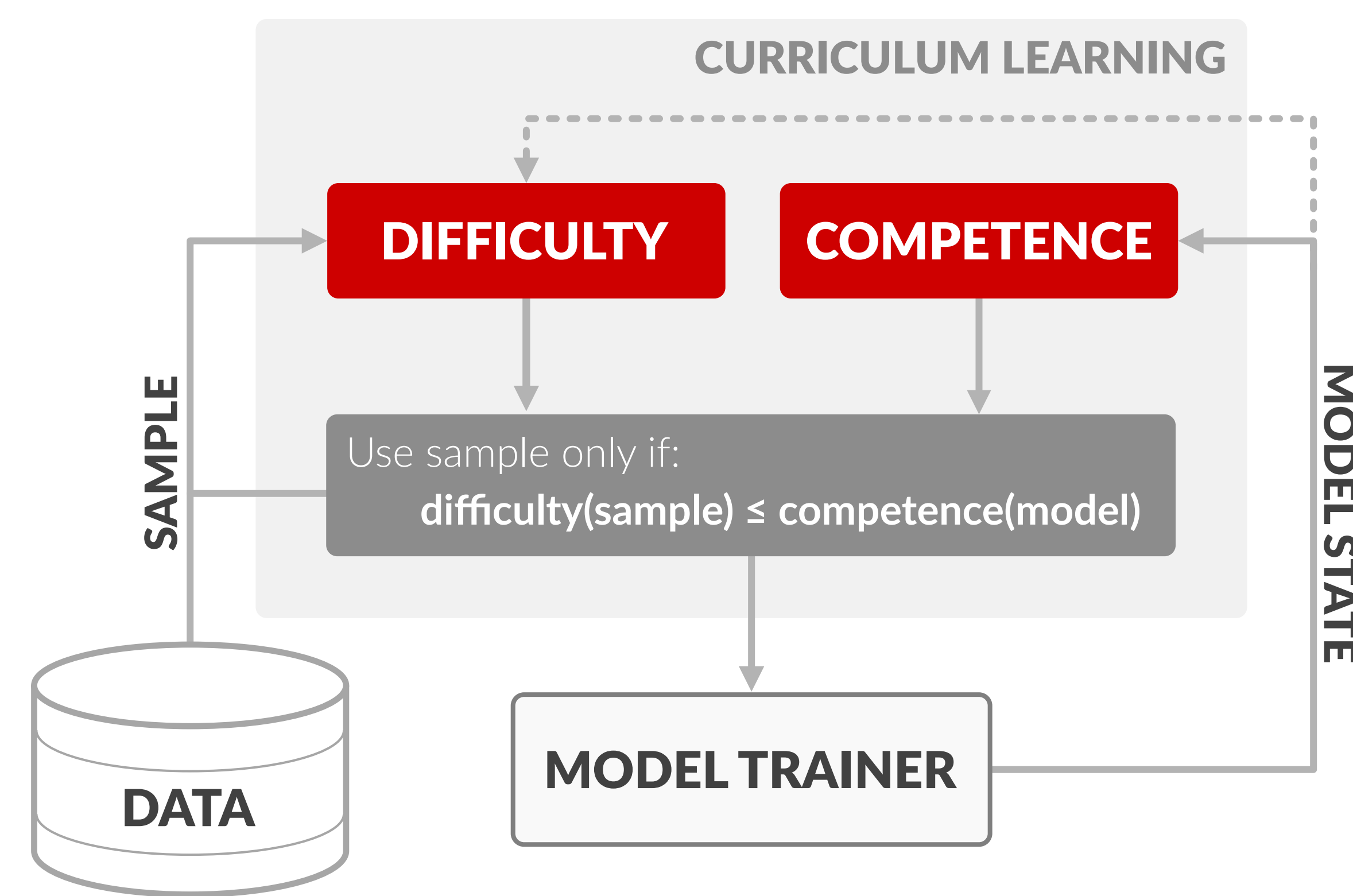


Sentence Length

$$d_{\text{length}}(s_i) \triangleq N_i \rightarrow \text{\#words in sentence}$$

Sentence Rarity

$$d_{\text{rarity}}(s_i) \triangleq - \sum_{k=1}^{N_i} \log \hat{p}(w_k^i)$$

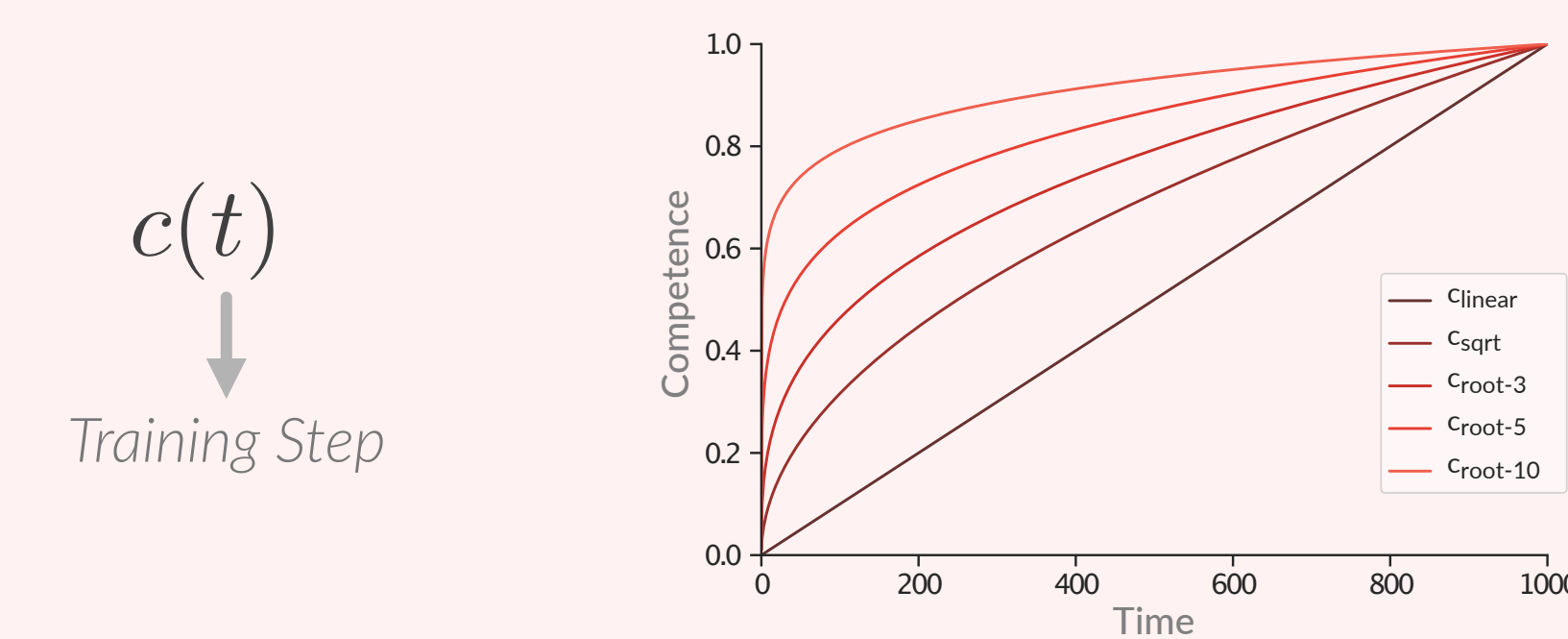


ALGORITHM

1. Compute the difficulty $d(s_i)$ for each sentence s_i .
2. Compute the cumulative density function (CDF), $\bar{d}(s_i) \in [0, 1]$ of the difficulties.
3. For training step $t = 1, \dots$:
 - i. Compute the model competence
 - ii. Sample a data batch uniformly from all examples such that: $\bar{d}(s_i) \leq c(t)$.
 - iii. Invoke the model trainer using the sampled batch.

COMPETENCE

Value between 0 and 1 that represents the progress of a learner during its training and can depend on the learner's state:



(e.g., validation set performance)

Linear Competence

New training examples are constantly being introduced during the training process with a constant rate $r = (1 - c_0)/T$, as a proportion of the total number of available training examples:

$$c_{\text{linear}}(t) \triangleq \min \left(1, t \frac{1 - c_0}{T} + c_0 \right)$$

time after which the learner is fully competent initial learning rate

Square Root Competence

Keep the rate in which new examples come in, inversely proportional to the training data size:

$$\frac{dc(t)}{dt} = \frac{P}{c(t)} \rightarrow c_{\text{sqrt}}(t) \triangleq \min \left(1, \sqrt{t \frac{1 - c_0^2}{T}} + c_0^2 \right)$$

Experiments

DATASETS

- IWSLT-15 (En→Vi)
- IWSLT-16 (Fr→En)
- WMT-16 (En→De)

MODELS

- RNN:
 - Bidirectional LSTM encoder
 - LSTM decoder
 - BPE vocabulary
- Transformer:
 - Base model of
 - BPE vocabulary Vaswani et al.

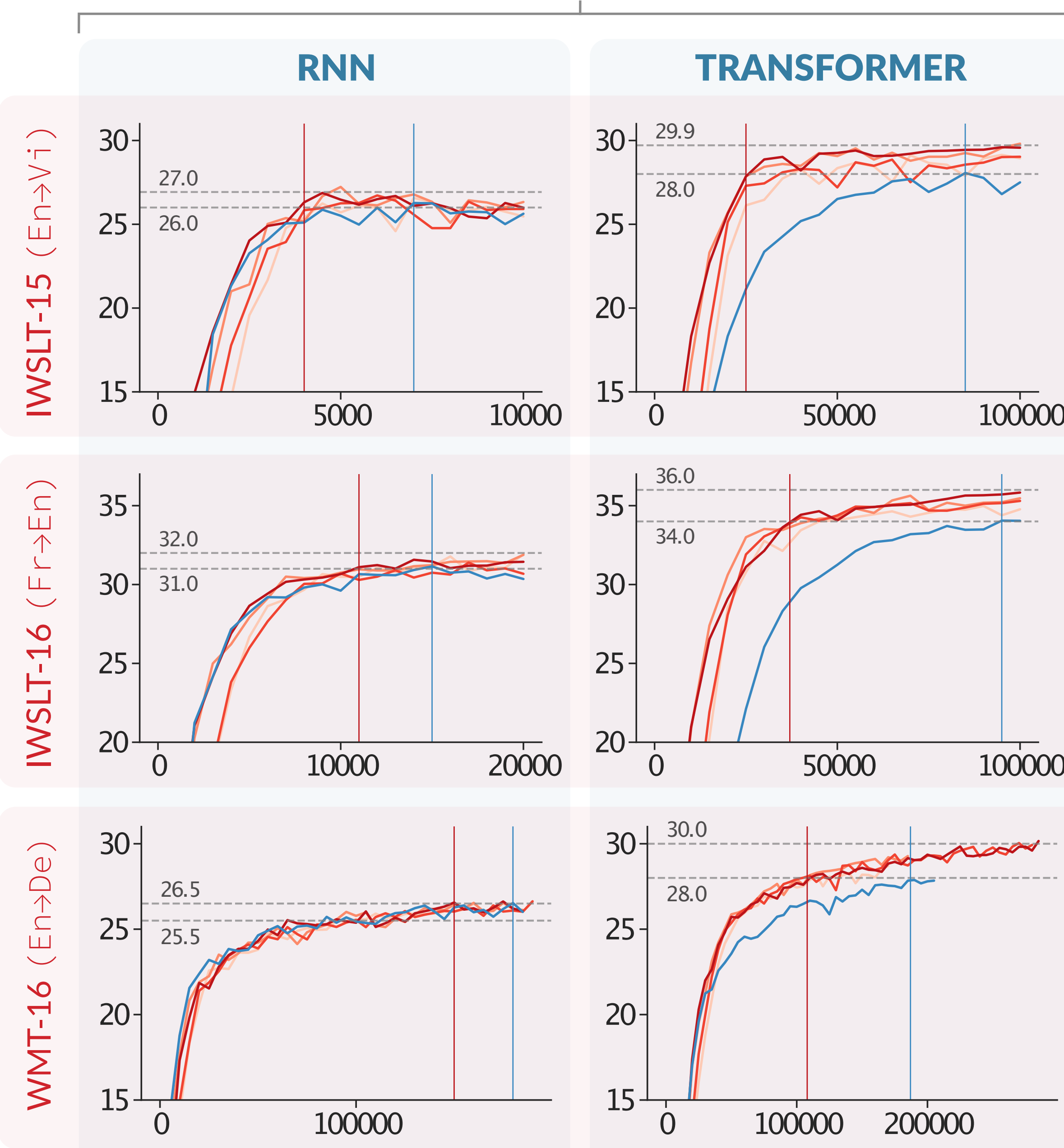
PARAMETERS

- Initial Competence: 0.01
- Curriculum Length: We train the baseline model without any curriculum, and compute the number of training steps it takes to reach ~90% of its final BLEU score.

CODE

Scala MT library to reproduce experiments:
<https://github.com/eaplatanios/symphony-nt>
TensorFlow Scala used for our experiments:
https://github.com/eaplatanios/tensorflow_scala

TRANSLATION BLEU SCORE



RELATIVE TIME TO BASELINE PERFORMANCE

