

Stochastic Transformer Networks with Linear Competing Units: Application to end-to-end SL Translation

Supplementary Material

Andreas Voskou^{*1}, Konstantinos P. Panousis¹, Dimitrios Kosmopoulos²,
Dimitris N. Metaxas³, and Sotirios Chatzis¹

¹Cyprus University of Technology, ²University of Patras, ³Rutgers University, New Jersey

1. Evidence Lower-Bound expression

The expression of the model Evidence Lower-Bound (ELBO) takes the form:

$$\mathcal{L}(q) = -\mathbb{E}_{q(\cdot)}[CE] - (\text{KL}[q(\boldsymbol{\xi})||p(\boldsymbol{\xi})] + \text{KL}[q(\boldsymbol{w})||p(\boldsymbol{w})]) \quad (1)$$

In this expression, CE represents the categorical cross-entropy loss obtained from the translation process viewed as a classification (selection of one-out-of-many) task. In the context of our model, this is a (stochastic) function of the model latent variables, $\boldsymbol{\xi}$ and \boldsymbol{w} . We have

$$\mathbb{E}_{q(\cdot)}[CE] = -\mathbb{E}_{q(\cdot)}[\log p(\mathcal{D}|\{\boldsymbol{w}, \boldsymbol{\xi}\})] \quad (2)$$

On the other hand, KL stands for the Kullback-Leibler divergences pertaining to the model latent variables, i.e. the winner unit indicator latent variables, $\boldsymbol{\xi}$, and the connection weights, \boldsymbol{w} . We have:

$$\text{KL}[q(\boldsymbol{\xi})||p(\boldsymbol{\xi})] = \sum_{\forall \boldsymbol{\xi}} \beta_{\boldsymbol{\xi}} \sum_{i=1}^U q(\xi_i) \log(q(\xi_i)U) \quad (3)$$

$$\text{KL}[q(\boldsymbol{w})||p(\boldsymbol{w})] = \sum_{\forall \boldsymbol{w}} \beta_w \left[\frac{\boldsymbol{\mu}^2 + \boldsymbol{\sigma}^2}{2} - \log \boldsymbol{\sigma} - \frac{1}{2} \right] \quad (4)$$

In these expressions, the factors $\beta_{\boldsymbol{\xi}}$ and β_w are heuristic hyperparameters needed to scale appropriately the KL contribution of each layer.

As we observe, the model ELBO entails posterior expectations over the network latent variables, $\boldsymbol{\xi}$ and \boldsymbol{w} . These cannot be computed analytically. Therefore, we approximate them via Monte-Carlo (MC) samples. To allow for low-variance ELBO gradients, we resort to the reparameterization trick.

Specifically, we express the samples of the connection weights, \boldsymbol{w} , in the form:

$$\boldsymbol{w} = \boldsymbol{\mu} + \boldsymbol{\sigma} \cdot \boldsymbol{z}, \quad \boldsymbol{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (5)$$

For the winner latent indicator variables, $\boldsymbol{\xi}$, the Gumbel-Softmax relaxation trick [1] yields the following expression for the drawn samples:

$$\boldsymbol{\xi} = \frac{\exp((\log \boldsymbol{\eta} + \boldsymbol{g})/T)}{\sum_{i=1}^U \exp((\log \eta_i + g_i)/T)} \quad (6)$$
$$\boldsymbol{g} = -\log(-\log \boldsymbol{z}), \quad \boldsymbol{z} \sim \text{U}(\mathbf{0}, \mathbf{1})$$

where $\boldsymbol{\eta} = q(\boldsymbol{\xi})$, and T is a positive temperature hyperparameter.

2. Effect of sample size S at inference time

As we explained in the main text, inference is performed through Bayesian averaging with a sample size of $S = 4$. It is useful to examine how this selection affects BLEU-4 scores. Thus, we repeat our experiments with the (2-2) configuration for different S values, and report the outcomes in Table 1. In addition, we exploit this opportunity to examine how the model behaves if we do not sample winners and connection weights; instead we pick the unit that yields the maximum posterior probability $q(\boldsymbol{\xi})$, and use the connection weight means, $\boldsymbol{\mu}$, to perform feedforward computations ("Deterministic" scenario).

We observe that the increased sample size indeed improves model stability and performance. More specifically, a large sample size $S = 16$ can slightly increase BLEU-4, while $S = 1$ produces a much lower performance than the proposed $S = 4$ size. However, large sample size also increases the model's inference time, and it may be impractical, especially for $S \gg 16$. Moreover, we can see

^{*}ai.voskou@edu.cut.ac.cy

Table 1. Sample size effect (BLEU-4 scores).

Sample	32 bit		Reduced	
	Dev	Test	Dev	Test
Deterministic (1)	22.99	22.83	22.21	22.31
1	23.45	23.05	22.66	23.03
4	23.23	23.64	23.09	23.52
16	23.34	23.95	23.14	23.39

that a deterministic version of our model is clearly inferior to stochastic operation, even if we draw just one sample, $S = 1$. This corroborates our theoretical intuitions regarding the doubly stochastic formulation of our model.

3. Computational Complexity.

We emphasize we do MC (simple drawing of i.i.d. samples; augmented with the reparameterization trick during training for low-variance gradients), and *not* Markov chain Monte-Carlo (MCMC), which of course is extremely inefficient. In table 2, we compare the computation times of our approach with the baseline. The differences in train-

Table 2. Computation Time on a single Quadro P5000 16GB.

	Baseline	Ours
Training (per batch)	0.2s	0.5s
Inference (single video)	0.04s	0.05s

ing time are due to the increased number of trainable parameters, namely the trainable variances of the Gaussian weights, and *not* due to Bayes. We emphasize that these trainable variances are used for proper weight compression, and are eventually discarded after training. Crucially, *convergence took almost the same number of epochs* with the baseline. Figure 1 illustrates both the learning curves and proves the latter statement and the smooth convergence of our network. The computational efficiency of our approach becomes even more apparent when we check inference times, which are comparable with the baseline Transformer when using $S = 4$ samples (that’s the best trade-off between accuracy and complexity).

Note also that, contrary to large Transformers that dominate NLP, our Transformer-type model is *small*, as we describe in the paper; just 2 layers! Crucially, the times reported in Table 2 were obtained on a *single* Quadro P5000 16GB GPU, which is inexpensive. Even more importantly, the compressed version of our network imposes a memory footprint of less than 1 GB. Thus, loading the trained model is feasible even on a commercial smartphone device.

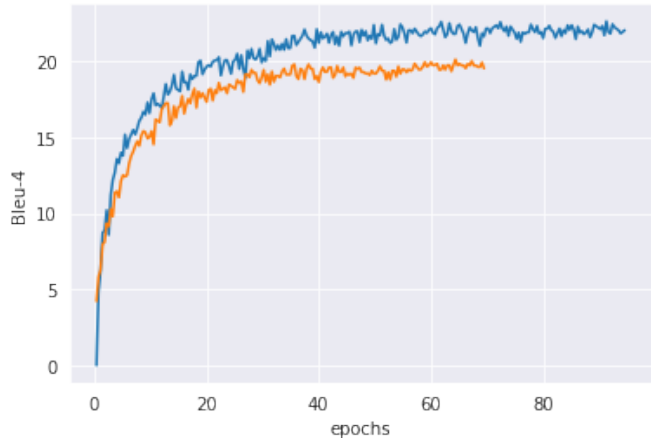


Figure 1. Convergence curves. Blue: proposed, Orange: baseline.

4. Translation examples

In Table 3, we present an extended sample of our outputs. We include both the German original and the corresponding translation in English.

References

[1] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax, 2017.

Table 3. Output sentences with the corresponding translations in English: Reference (R), single model (S), and ensemble (E).

<p>R: vor allem im süden können die gewitter unwetterartig ausfallen . (<i>especially in the south, thunderstorms can be severe.</i>)</p> <p>S: vor allem im süden und süden muss mit sturmböen gerechnet werden . (<i>especially in the south and the south, gusts of wind must be expected.</i>)</p> <p>E: vor allem im süden und süden muss mit unwetterartigen gewittern gerechnet werden . (<i>especially in the south and south of the country, thunderstorms are to be expected.</i>)</p>
<p>R: am donnerstag vor allem in der nordwesthälfte hier und da sturmböen . (<i>on thursday, squalls here and there, especially in the northwestern half.</i>)</p> <p>S: am donnerstag vor allem im norden und westen teilweise kräftig . (<i>on thursday especially in the north and west partly strong.</i>)</p> <p>E: am donnerstag vor allem im norden und westen teilweise frischer wind . (<i>on thursday partly fresh wind, especially in the north and west.</i>)</p>
<p>R: morgen bei dauernebel im süden nur zwei am niederrhein bis elf grad . (<i>tomorrow with permanent fog in the south only two on the lower Rhine up to eleven degrees.</i>)</p> <p>S: morgen im süden bis zweiundzwanzig am oberrhein bis elf grad . (<i>tomorrow in the south up to twenty-two degrees on the upper rhine up to eleven degrees .</i>)</p> <p>E: morgen im süden noch zwei bis einundzwanzig am rhein bis elf grad . (<i>tomorrow in the south still two to twenty-one on the rhine to eleven degrees.</i>)</p>
<p>R: im erzgebirge morgen dreizehn sonst fünfzehn bis zweiundzwanzig grad . (<i>in the ore mountains tomorrow thirteen or fifteen to twenty-two degrees .</i>)</p> <p>S: an der saale morgen dreizehn in der eifel bis zweiundzwanzig grad . (<i>on the Saale tomorrow thirteen in the Eifel to twenty-two degrees .</i>)</p> <p>E: an der see morgen dreizehn sonst zwei bis zweiundzwanzig grad . (<i>at the lake tomorrow thirteen otherwise two to twenty-two degrees .</i>)</p>
<p>R: in der neuen woche kühlt es dann bei wechselhaftem wetter deutlich ab . (<i>in the new week it cools down considerably with changeable weather.</i>)</p> <p>S: in der neuen woche wechselhaft und deutlich kühler . (<i>in the new week changeable and clearly cooler.</i>)</p> <p>E: in der neuen woche unbeständig und noch kühler . (<i>in the new week unstable and even colder.</i>)</p>
<p>R: und so warm mit zwanzig grad wird es nicht mehr in den nächsten tagen . (<i>and it won't be that warm with twenty degrees in the next few days.</i>)</p> <p>S: und die zwanzig geht es in den nächsten tagen . (<i>and the twenty it goes in the next few days.</i>)</p> <p>E: und die zwanzig grad in den nächsten tagen . (<i>and the twenty degrees in the next few days.</i>)</p>
<p>R: mit föhnunterstützung klettern die temperaturen am alpenrand bis zehn grad . (<i>with foehn support, temperatures climb to ten degrees on the alpine fringe.</i>)</p> <p>S: es ist ein bisschen kühler temperaturen so ein bisschen verhalten zehn bis zehn grad . (<i>it is a bit cooler temperatures so a bit restrained ten to ten degrees.</i>)</p> <p>E: es kann ein bisschen geben bei temperaturen am alpenrand da sinken die temperaturen bis zehn grad . (<i>there can be a bit of a dip in temperatures at the edge of the alps, where temperatures drop to ten degrees.</i>)</p>
<p>R: mit der leicht kühleren luft dann bis vierundzwanzig grad in der nordhälfte . (<i>with the slightly cooler air then up to twenty-four degrees in the north half.</i>)</p> <p>S: in diesem bild von vierundzwanzig bis vierundzwanzig grad in der lausitz . (<i>in this picture from twenty-four to twenty-four degrees in Lusatia .</i>)</p> <p>E: in der zweiten nachthälfte die vierundzwanzig bis vierundzwanzig grad in der lausitz . (<i>in the second half of the night the twenty-four to twenty-four degrees in Lusatia .</i>)</p>
<p>R: auch in den folgenden tagen ändert sich an diesem wechselhaften wetter wenig . (<i>even in the following days there is little change in this changeable weather.</i>)</p> <p>S: in den folgenden tagen setzt sich das wechselhafte witterung fort . (<i>in the following days the changeable weather continues.</i>)</p> <p>E: in den folgenden tagen setzt sich das wechselhafte wetter fort . (<i>the weather continues to be changeable in the following days.</i>)</p>