

Federal State Autonomous Educational Institution for Higher Education  
National Research University Higher School of Economics

Faculty of Computer Science  
Educational Program  
Applied Mathematics and Information Science

TERM PAPER  
RESEARCH PROJECT  
"UNCERTAINTY ESTIMATION FOR MACHINE TRANSLATION"

Prepared by the student of group 171, 3rd year of study  
Kuznetsov Dmitriy Sergeevich

Supervisor  
research fellow Lobacheva Ekaterina Maksimovna

Moscow 2020

# Contents

<b>1</b>	<b>Annotation</b>	<b>2</b>
<b>2</b>	<b>Introduction</b>	<b>3</b>
<b>3</b>	<b>Machine Translation</b>	<b>4</b>
3.1	Definition . . . . .	4
<b>4</b>	<b>Beam Search</b>	<b>5</b>
4.1	Method definition . . . . .	5
4.2	Large Beam Size Problem . . . . .	6
<b>5</b>	<b>Uncertainty</b>	<b>8</b>
5.1	Definition . . . . .	8
5.2	Uncertainty Estimation for Structured Predictions . . . . .	8
5.3	Analyzing Uncertainty in NMT . . . . .	9
<b>6</b>	<b>Summary</b>	<b>10</b>
<b>7</b>	<b>Baseline approach</b>	<b>11</b>
<b>8</b>	<b>Misclassification detection</b>	<b>12</b>
8.1	Experiment correctness . . . . .	17
8.2	Using on practice . . . . .	19
8.3	Beam search . . . . .	19
<b>9</b>	<b>In-ensemble uncertainty estimate</b>	<b>20</b>
9.1	Probability distribution calibration . . . . .	23
9.2	Calibration BLEU . . . . .	26
<b>10</b>	<b>Conclusion</b>	<b>28</b>

# 1 Annotation

Machine learning models uncertainty estimation is an important task in different areas. Uncertainty estimation helps to improve system quality, marking inconsistent parts of models behavior. For such task as ASR or image classification uncertainty estimation is well investigated, but in the machine translation field uncertainty estimation is poorly explored. Machine translation data uncertainty estimation task has been already considered. In this work we investigate total and model uncertainty estimation in neural machine translation for models ensemble. We explore correlation between uncertainty and error in model's predictions.

***Index terms*** - *machine translation, uncertainty estimation, beam search, larger beam size problem, model uncertainty, data uncertainty*

***Abstract in russian*** - Оценка неопределенностей моделей машинного обучения - важная задача в различных областях. Оценка неопределенностей помогает улучшить качество программной системы, помечая несогласованные аспекты поведения моделей. Для таких задач, как распознавание речи или классификация изображений, оценка неопределенностей хорошо изучена, но в задаче нейро-машинного перевода оценка неопределенностей исследована плохо. На данный момент существуют работы исследующие оценку неопределенностей данных в задаче нейро-машинного перевода. В данной работе мы исследуем полную и модельную оценку неопределенностей в задаче нейро-машинного перевода для ансамбля моделей. Мы изучим корреляцию между неопределенностью и ошибкой в предсказаниях моделей.

## 2 Introduction

There are two fundamental approaches in machine translation: statistical and *neural machine translation* (NMT, further). Neural machine translation models achieve state-of-the-art results. As in any field, there are a lot of actual unresolved problems in machine translation. In particular, in neural machine translation quality of prediction degradation appears with a hypothesis search width increase.

Neural models have significant problems with confidence in token probability. *Uncertainty estimation* is a way to measure models' confidence level. While we can estimate model's confidence, it is possible to re-organise final prediction according to uncertainty estimation.

There are several types of uncertainties. Some of them are caused by certain artifacts in the data, while others are caused by the model's inability to understand the proposed data. In any case, the uncertainty estimates indicate that the model is doubtful about making certain predictions. Moreover, we can use knowledge about uncertainty to improve prediction quality, e.g. as we introduce later it is possible to calibrate output softmax distribution to avoid misclassifications in models' predictions.

In this paper, we study the correlation between machine translation model errors and uncertainties. We empirically prove the fact, uncertainty estimates can be a misclassification indicator. We explore properties of total, data and model uncertainty estimates and its connection with correct token and error token probabilities.

Moreover, we propose an approach to using uncertainty estimation as a potential way to improve the quality of the model or architecture. We present token probability calibration method based on inensemble models' probability variance. It is not improve target quality, but we investigate significant observations, which could be usefull in further works. We explore, how to construct such calibration function, and its influence on token probability distribution.

## 3 Machine Translation

### 3.1 Definition

The task of machine translation consists of providing a sentence equivalent in meaning in the target language for the original sentence in the source language. Let  $e_1, \dots, e_n$  be a sequence of an input sentence tokens' vector representation in the source language. The machine translation model is required to build a sequence of tokens  $f_1, \dots, f_m$  in the target language. Moreover, in a target language, the sequence must have the original meaning.

Most modern neuromachine translation models follow the paradigm [Sutskever et al. \(2014\)](#) (seq2seq). Seq2seq models often consist of two networks or groups of networks. One network that processes input tokens is called *encoder*, the second network that builds output predictions is called *decoder*. The encoder encodes an input sequence into some hidden representations  $h_1, \dots, h_p$ , which are fed to the decoder as hidden states before building the output prediction. The decoder output is a vector set  $\{(p_{i1}, \dots, p_{id})\}_{i=1}^m$ , where  $d$  - target language dictionary power,  $m$  - maximum possible prediction sequence length,  $p_{ij}$  - the model's confidence that the  $i$ th token of the predicted sentence is the  $j$ th token from the target language dictionary. As a result, the translation model builds a "probability distribution" in the space of the cartesian product of the target dictionary. On practice, it is hard to choose final hypothesis using joint tokens' probability distributions (too much variants). That is why we construct final hypothesis token by token using conditional probability. Let us consider *autoregressive* model:

$$P(\mathbf{y}|x, \theta) = P(y_1|x, \theta) \prod_{i=2}^L P(y_i|y_{<i}, \mathbf{x}, \theta) \quad (1)$$

where  $\mathbf{x}$  - input sequence,  $\mathbf{y}$  - model's final translation,  $\theta$  - model's training parameters,  $y_{<i} = (y_1, \dots, y_{i-1})$ .

Autoregressive model constructs final hypothesis step by step. For each token  $y_i$  model considers token probability distribution and its cumulative distribution

and choose next translation token according to some rule.

The simplest way to build a prediction (final translation) is to select the token with the highest probability  $f_i := \operatorname{argmax} \{p_{ij}\}_{j=1}^d$  at each step. So  $f_1, \dots, f_m$  is a sequence of model hypotheses' tokens. However, for an optimal prediction  $f_1^*, \dots, f_m^*$  it is not always true that in terms of model confidence  $f_i^* = \operatorname{argmax} \{p_{ij}\}_{j=1}^d$  (it is true only if tokens are independent). This phenomenon is also related to the beam problem, which we discuss in more detail in the following chapters.

The most common quality measure in a machine translation problem is BLEU.

Let the model for some pair of sentences be  $\mathbf{x}, \mathbf{r}$  built a translation of  $\mathbf{y}$ .  $\mathbf{x}$  - input sentence in a source language and  $\mathbf{r}$  - its ground truth translation.

$$BLEU_K = \min \left( 1, \frac{|\mathbf{y}|}{|\mathbf{r}|} \right) \left( \prod_{i=1}^K precision_i \right)^{\frac{1}{K}} \quad (2)$$

where,  $precision_i$  means *precision* calculated for all i-grams.

## 4 Beam Search

### 4.1 Method definition

Beam search is a method to obtain a more effective solution for building predictions. Instead of selecting most likely token at each step, we require the model to store  $b$  the most likely predictions' prefixes  $y_{b1}, \dots, y_{bt}$ . Using this approach, we do not choose the locally optimal token, but in general the optimal prefix, which increases the chance of building a correct prediction. It is important to notice, beam search still does not guarantee, that final translation is an optimal in  $\mathcal{D}^L$  space, but it broadens search space comparing to greedy-search.

Let us define  $k$ -long prediction prefix:

$$y_{<k}^i := (y_1^i, \dots, y_{k-1}^i), \quad \text{where } i - \text{beam search branch index} \quad (3)$$

Suggest we have already finished  $k$  iterations of beam search. At that moment we find *beam* most likely  $k$ -long prediction prefixes:

$$y_{<k+1}^1, \dots, y_{<k+1}^{beam}$$

Beam Search makes  $k + 1$  step as follow:

- 1 For all  $i \in [1, beam]$  beam search consider cumulative probabilities

$P_i = \{ p(y_{k+1}^i | y_{<k+1}^i, \mathbf{x}, \theta) p(y_{<k+1}^i | \mathbf{x}, \theta) \mid y_{k+1}^i \in \mathcal{D} \}$  and corresponding  $y_{k+1}^i$  tokens.

- 2 Let  $P = \bigcup_{i=1}^{beam} P_i$

- 3 Beam search chooses *beam* most highest probabilies in  $P$  and its corresponding prefixes  $y_{<k+1}^i$  and its corresponding continuations  $y_{k+1}^i$ .

- 4 So we have  $y_{<k+2}^1, \dots, y_{<k+2}^{beam}$   $k + 1$ -long prefixes on  $k + 1$  step.

Finally, we complete  $L$  autoregressive iterations of Beam Search. We have  $y_{<L+1}^1, \dots, y_{<L+1}^{beam}$  prefixes. Final translation is most likely  $L + 1$ -long prefix among them.

In summary, this method stores *beam* of the "most likely" prefixes at each iteration, in contrast to the naive approach, which greedily selects one locally optimal hypothesis at each iteration.

## 4.2 Large Beam Size Problem

In NMT tasks often there is a problem, that starting from a certain threshold value, the quality metric decreases with the growth of the beam size. It calls *large beam size problem*.

[Murray et al. \(2018\)](#) raises two issues: the large beam size problem and the tendency of NMT models to make short predictions. According to the authors, solving the short prediction problem involves solving the beam problem.

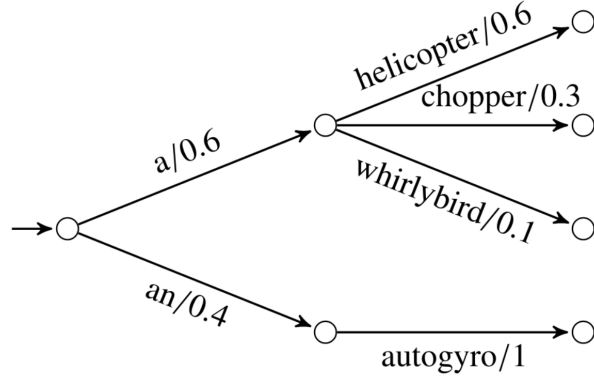


Figure 1 – Label bias causes this toy word-by-word translation model to translate French helicopter incorrectly. This figure is taken from [4].

As a demonstration of the reasons listed in this article that cause the beam problem, consider the following example from the article.

On the Fig. 1. we can see a tree of Beam Search predictions for the word *un hélicoptère*. Let the beam size be equal to two and on the first iteration we save two hypotheses: *a* and *an*. Note that all 4 following hypotheses can be potential translation options, but *a helicopter* in this example is a correct translation. However, *autogyro* is the only one continuation for prefix *an*, therefore the model confidence for this suffix is 1, and as a result, the model confidence to predict *an autogyro* is equal to 0.4. At the same time, the probability mass for the *a* continuation suffixes is divided among three words: *helicopter*, *chopper*, *whirlybird*. The confidence for the correct translation of *a helicopter* from the model’s point of view is 0.36. As a result, the model gives a incorrect prediction. The larger beam size we take, the higher probability we get to meet this type of situations, hence the loss of quality.

Note that this example demonstrates the fact that if a certain prefix involves a low-entropy suffix, the model tends to ignore the second input token and predict with high confidence one of several highly probable low-entropy suffixes. As a result, the low entropy of the suffix leads to overestimation of this prediction branch and poor quality. In addition, there is a high probability of getting short predictions due to low entropy, since the model begins to ignore some tokens in favor of high confidence of low-entropy suffixes. For this reason, the authors of



the article consider the short prediction problem and the beam problem together.

## 5 Uncertainty

### 5.1 Definition

In the [Malinin et al. \(2018\)](#) authors give following uncertainty definition. Uncertainty in machine learning is usually understood as a measure of the non-degeneration of the distribution on the model predictions at a fixed input. Less formally, if we were able to estimate uncertainty, we could tell how much we can trust the model’s predictions.

There are several types of uncertainties: data uncertainty, model uncertainty, and distributional uncertainty. *Data uncertainty (aleatoric uncertainty)* is a type of uncertainty caused by the nature of the data, including noise, class balance, etc. *Model uncertainty (epistemic uncertainty)* measures uncertainty caused by a model specificity, how good a model understands given data. *Distributional uncertainty* measures dissimilarity between training and test data, caused by unknown or strange test examples.

If we want to estimate uncertainty, it is important for the model to evaluate all types of uncertainty, because they are caused by different reasons and have different negative effects.

### 5.2 Uncertainty Estimation for Structured Predictions

In this section, we discuss various approaches to estimating uncertainty for structured predictions and some its applications, proposed by [Malinin et al. \(2020\)](#). In this section, we consider problems in which solutions are described by the probabilistic model (1).

The author identifies two main approaches to estimating uncertainties for structured predictions: sequence-level and token-level. In the case of sequence-level, the uncertainty is estimated in the space of final predictions (the Cartesian

product of the dictionary). In the case of token-level, the uncertainty is estimated at the level of individual sequence tokens (the dictionary space). Estimating uncertainties at different levels allows us to solve different types of problems.

The main tool for estimating uncertainty in the article is the use of models ensemble from a certain parametric family of models  $\mathcal{F}(\theta)$ . If many different models from the same family vote for the same verdict, it is highly probable that we can trust this verdict. In the opposite situation, if the models are not consistent in their predictions, this fact can become an indicator of high uncertainty. In particular, it may indicate that the predicting token is out of the general pattern of the data that the models were able to fit in their generality.

Authors suggest to estimate uncertainty using following realation:

$$\underbrace{\mathcal{I}[y, \theta | \mathbf{x}, \mathcal{D}]}_{\text{Knowledge Uncertainty}} = \underbrace{\mathcal{H}[P(y | \mathbf{x}, \mathcal{D})]}_{\text{Total Uncertainty}} - \underbrace{\mathbb{E}_{q(\theta)}[\mathcal{H}[P(y | \mathbf{x}, \theta)]]}_{\text{Data Uncertainty}} \quad (4)$$

Here  $\mathcal{H}$  is Shannon entropy and  $\mathcal{I}$  - mutual information.

Knowledge uncertainty can also be evaluated using a different metric suggested in the article, which correlates with the one already considered:

$$\mathcal{K}[y, \theta] = \mathbb{E}_{q(\theta)q(\hat{\theta})}[\text{KL}[P(y | \mathbf{x}, \theta) || P(y | \mathbf{x}, \hat{\theta})]] \quad (5)$$

Authors demonstrate the effectiveness of uncertainty estimation of described methods at various estimation levels for the *Artificial Speach Recognition* (ASR) problem. In particular, experiments: out of distribution detection and misclassification detection. However, this article does not deal with token-level uncertainty estimation for a machine translation problem. In our work we conduct experiments with token-level uncertainty estimation for NMT.

### 5.3 Analyzing Uncertainty in NMT

*Data uncertainty* was also investigated by [Ott et al. \(2018\)](#).

Authors consider two types of data uncertainty. *Intrinsic uncertainty* is the re-

sult of existence of several semantically equivalent translations of the same source sentence, for instance there are several ways to express the same meaning. Situations when target language is more complex than a source language, can also cause a high uncertainty. In such situations it could be impossible to earn additional information, such as gender or tense, which is necessary to build a correct prediction in a target language. Web data crawling for train corpus construction can cause *extrinsic uncertainty*. For example, authors notice that at least 1% sentences in WMT datasets is just a copy of a source sentences.

One of the main problems raised in the article is an effect of *uncertainty* on the degradation of quality with the growth of beam size. Authors notice, that sentences copies (which cause extrinsic uncertainty) are overrepresented in the output of beam search. The larger beam size one takes, the more copies appear in a beam search tree.

Authors make a conclusion, high copies ratio correlate with BLEU metric degradation. As we mentioned above, larger beam size leads to high probability to face with copies in the beam search output. So larger beam size leads to BLEU degradation.

So using extrinsic uncertainty estimation can improve model quality. We can clean up dataset, decreasing uncertainty, which leads to BLEU increasing.

## 6 Summary

In this paper, we study the correlation between uncertainty and translation errors. In particular, we want to find out how to improve the quality of translation by adding some knowledge about the uncertainty of prediction to the model. We consider token-level estimation approach, because of this approach for the machine translation task is poorly explored. We also analyse the correlation between the uncertainty estimation and the larger beam size problem.

## 7 Baseline approach

Let  $X = (x^1, \dots, x^N)$ ,  $R = (r^1, \dots, r^N)$  a parallel corpus of sentences, where  $X$  corresponds to the source language, and  $R$  corresponds to the target language (the one we are translating to). The final translation of a certain sentence  $x$  is based on the following probabilistic model:

$$P(\mathbf{y}|x, \theta) = P(y_1|x, \theta) \prod_{i=2}^L P(y_i|y_{<i}, x, \theta)$$

The purpose of this work is to use uncertainty estimation (total, knowledge) to determine situations when the model makes an error in predicting the next token  $y_i$ . We look at several approaches and ways to use uncertainty estimates.

As a baseline model we use the transformer [Vaswani et al. \(2017\)](#), implemented in the [Fairseq \(Facebook\)](#). In all experiments, we consider an ensemble of 5 transformers. The final probability model of the ensemble has a form:

$$P(y_i|y_{<i}, x, \mathcal{D}) = \frac{1}{K} \sum_{j=1}^K P(y_i|y_{<i}, x, \theta_j)$$

In our case,  $K = 5$ . Each model in the ensemble is trained on shared data  $X, R$ , but from different initialization points.

We use [IWSLT'14 De-En \(\)](#) dataset with standard train-val-test split for transformer models. We measure translation quality with *BLEU4*.

We use the following hyperparameters for training each model in the ensemble:

1 *clip\_norm* = 0

2 *learning\_rate* =  $5 * 10^{-4}$

3 *label\_smoothing* = 0.1

4 *epochs* = 75

5 *lenpen* = 0 (We want to know how uncertainty estimation affects Bleu degradation by beam size)

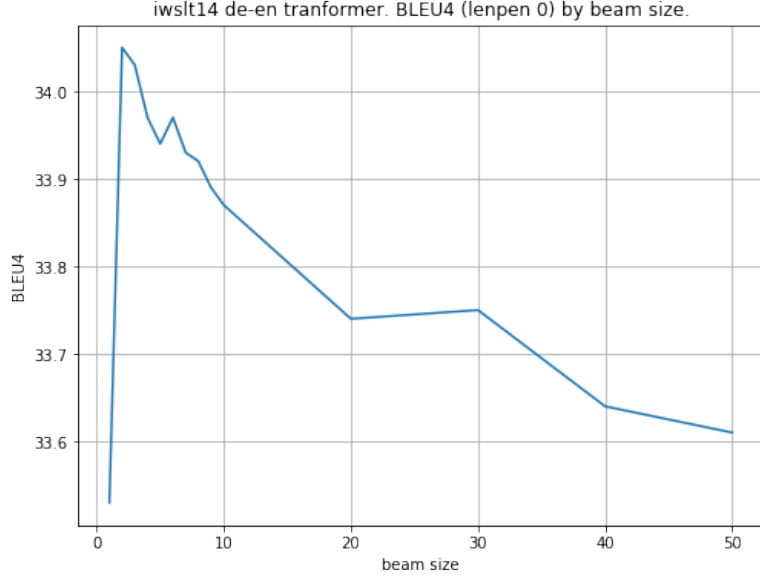


Figure 2 – BLEU on test set for transformer model by beam size

Single baseline model BLEU can be seen on Fig. 2. With the growth of beam size BLEU degrades. As we found out earlier, this can be caused by overestimating or underestimating the probabilities of certain tokens. In this paper, we also analyse the correlation of this problem with the uncertainty.

## 8 Misclassification detection

In [Malinin et al. \(2020\)](#), an attempt to detect falsely predicted tokens on the ASR problem is considered. Unfortunately, the author did not conduct a similar experiment for NMT due to the complexity of defining the correct target metric in machine translation.

In our work, we decided to explicitly check whether some uncertainty estimates can detect false predictions of the next token.

The structure of the experiment is as follows. We consider the corpus of sentences  $X, R$ . Let  $y = (y^1, \dots, y^N)$  - corresponding translations constructed by the ensemble. Let  $r_i, y_i$  be  $i$ th true translation and model translation tokens, respectively. Enter the following indicator of the  $i$ th token correctness in translation:

$$correct(y^i, \mathbf{r}) := [y^i \in \mathbf{r}] \quad (6)$$

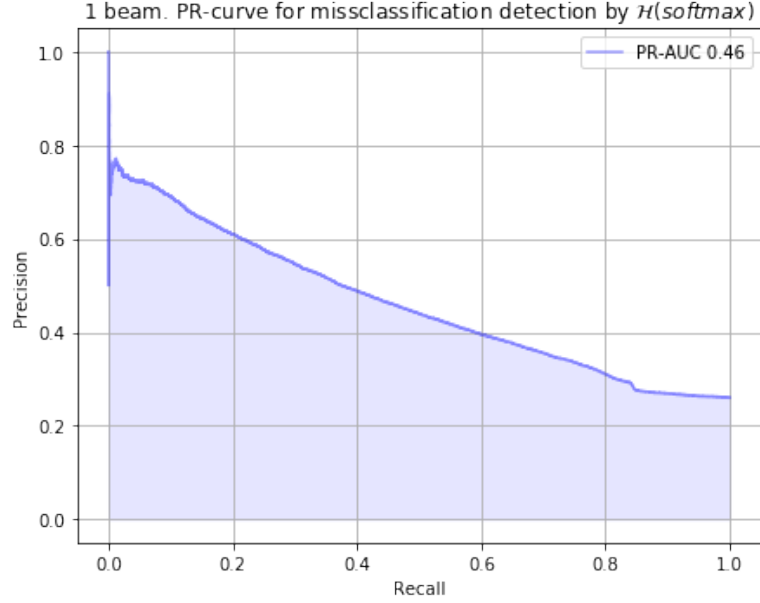


Figure 3 – PR-curve for missclassification detection by TU

This indicator value is our target function, which we want to predict based on the uncertainty estimation. Authors introduce the *total uncertainty estimation*:

$$TU(y_i) := \mathcal{H}[P(y_i|y_{<i}, \mathcal{D})] \quad (7)$$

Here, the Shannon entropy is considered as  $\mathcal{H}$ . This metric does not separate types of uncertainties that we have discussed earlier, it is an indicator for any manifestation. Metric show how much the probability mass of the ensemble is sparse by dictionary tokens.

Then for each  $y_i^j \in H$  token, we have the target functions  $correct(y_i^j, r^j)$  and uncertainty estimates  $TU(y_i^j)$ . We scale the uncertainty estimates to  $[0, 1]$ . So we get some model that predicts error probability of a particular token prediction.

Results of the experiment can be seen in Fig. 3. If we try to maximize the F1-score, results turned out to be worse than the random predictor. However, if we are intresting only in a small percentage of errors, we can get a relatively good detection accuracy.

Obtained results are consistent with the thesis put forward by article authors. It would seem that we have fixed easy correctness indicator (any token that is not in the translation at all is marked as an error), but the uncertainty estimation

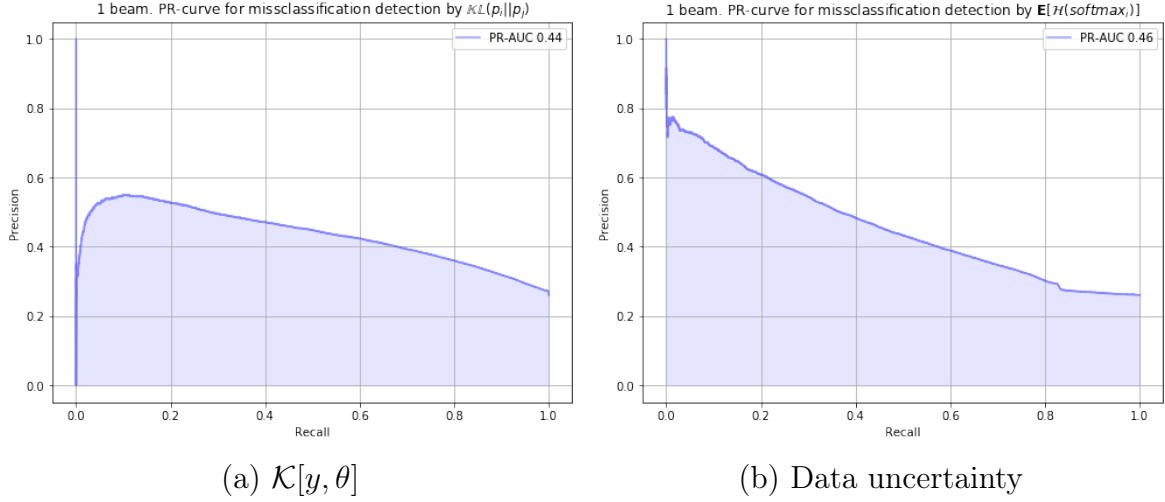


Figure 4 – PR-curve for misclassification detection

works worse than a random prediction. The problem lies in the fact that if the translation is a close meaning synonym, in which the ensemble models are strongly confident, then from the point of view of our indicator - this is an error, but the uncertainty estimation does not detect any anomalies. As noted by [Malinin et al. \(2020\)](#), the solution to the problem is a human manual data markup to identify more meaningful targets in this experiment.

We may be able to use other ratings or similarity indicators to achieve higher quality.

[Malinin et al. \(2020\)](#) consider another way of uncertainty estimation:

$$\mathcal{K}[y, \theta] = \mathbb{E}_{q(\theta)q(\hat{\theta})}[\text{KL}[P(y|\mathbf{x}, \theta) || P(y|\mathbf{x}, \hat{\theta})]]$$

Recall that this metric estimates knowledge uncertainty (model uncertainty). Let us take a look at experiment results for this case (Fig. 4a). As we can see  $\mathbb{KL}$  solves misclassification problem worse than previous metric. This results conform with the result that authors get in their paper for ASR model (Knowledge uncertainty metrics estimate worse than total uncertainty estimations). Nevertheless,  $\mathbb{KL}$  could be effective, if we define more resonable correctness indicator, as a total uncertainty estimation could be.

One more uncertainty estimation from the paper is data uncertainty estima-

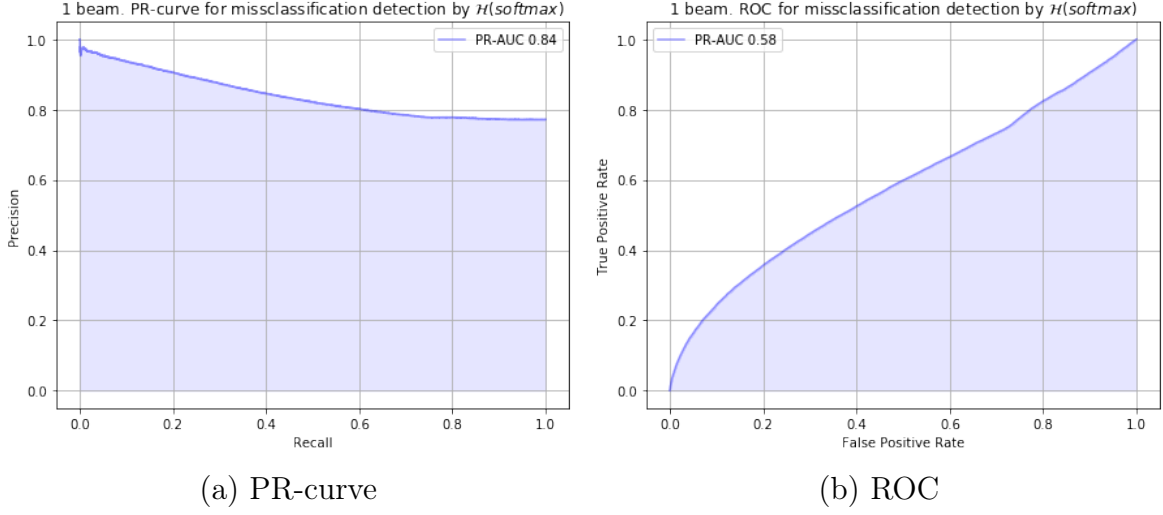


Figure 5 – Misclassification detection by TU.  $\delta(y^i, \mathbf{r})$

tion:

$$\mathbb{E}_{q(\theta)}[\mathcal{H}[P(y|\mathbf{x}, \theta)]]$$

Experiment results, using data uncertainty estimation, are presented in Fig. 4b. Obtained score is as bad as in TU case (as previous one, it conforms with Malinin et al. (2020) results). As we can see, data uncertainty estimation gains more score, than model uncertainty estimation. That is why we can suggest that, if considered metrics is reasonable, so uncertainty in data influences on model prediction quality more than uncertainty in models generalization ability in ensemble.

Total uncertainty, data uncertainty and knowledge uncertainty estimates low results caused by same reasons, that we discussed above.

Let us now conduct similar experiments for the new correctness indicator:

$$\delta(y^i, \mathbf{r}) := [y^i = r^i] \quad (8)$$

And consider results in Fig. 5a. This time there are good results. However, they are not caused by the fact that we have selected a more reasonable correctness indicator, but by the fact that this indicator defines the concept of error more strictly. Now it is more likely to meet a misclassification and as a result, the uncertainty estimates have become easier to deal with this. Therefore, these results are not very reasonable. This is easy to see if we consider the ROC in Fig.



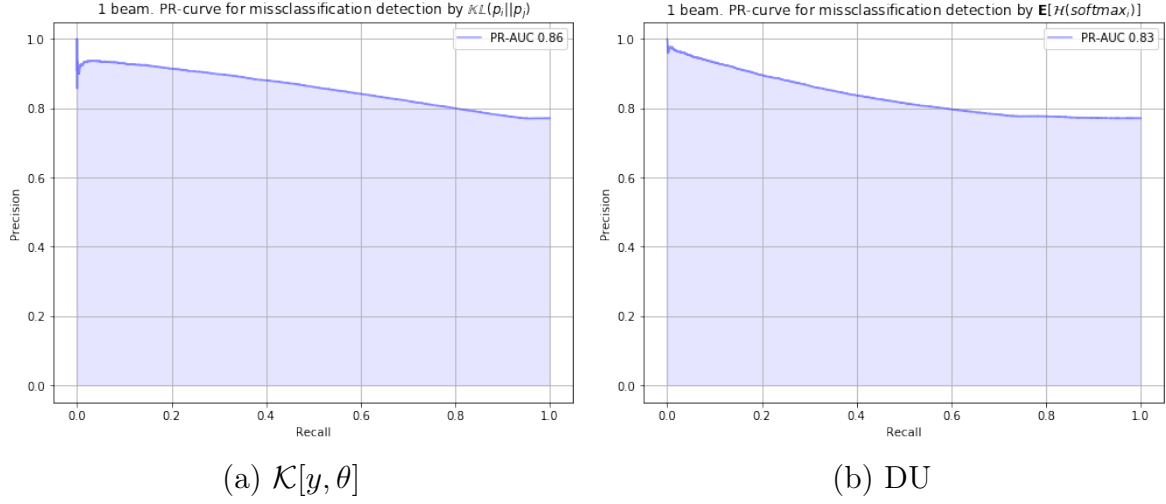


Figure 6 – Misclassification detection.  $\delta(y^i, \mathbf{r})$

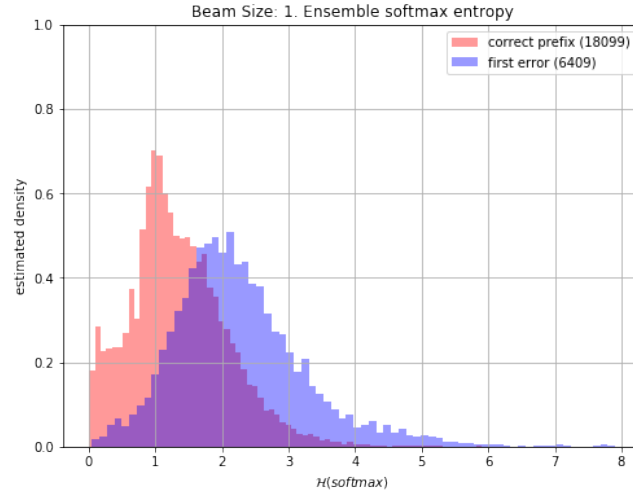


Figure 7 – Ensemble softmax Shannon entropy

5b. Our empirical model does not do much better than a random algorithm. It is important to note that after we strengthened the correctness condition, the percentage of positive class increased significantly (0.76 against 0.26 in the previous statement). In this statement, we can trust ROC, since the positive class prevails and as a result, the balance of classes does not worsen ROC-score.

As for previous correctness indicator, we consider data uncertainty estimate and KL estimate on Fig. 6.

In all considered cases results are unreasonable. We can not trust metrics, while we use groundless correctness indicator.

## 8.1 Experiment correctness

An important question arises: why a classification error can correlate with the uncertainty estimate TU.

Let us try to estimate a density of the softmax entropy distribution for two cases: we correctly predict the token and we made the first error in the sentence (we predicted the entire prefix before this token correctly). In the other words, for each such case, we calculate the softmax entropy and construct empirical density estimates (normalized histograms). The results can be seen in Fig. 7.

As we can see, the densities in both cases are distributed almost normally, unimodally. Moreover, the entropy distribution mode of correct tokens lies strictly to the left of the distribution mode for the first error, i.e. most of the probability mass of correct tokens has a lower softmax entropy (the TU metric).

How can we interpret this? The higher ans entropy of the distribution - the more uniform distribution is, i.e. more and more competing hypotheses appear, i.e. our ensemble identifies more and more equally probable variants, namely, less confident in its answer. Conversely, if the entropy is very low, then the distribution is close to singular, i.e. one hypothesis has about 1. confidence from the point of view of the ensemble.

Thus, if the TU is large enough, it can be an indicator that we are not very confident in the token we are predicting, which allows us to estimate uncertainty in prediction.

Further, let us collect same distribution but for:

$$\mathcal{K}[y, \theta] = \mathbb{E}_{q(\theta)q(\hat{\theta})}[\text{KL}[P(y|\mathbf{x}, \theta) || P(y|\mathbf{x}, \hat{\theta})]]$$

KL-divergence measures difference between two distributions (more - further), this estimate is always non-negative. While we take expectation of estimate, we try to average difference between softmaxes among all possible models pairs in ensemble. We can interpret  $\mathcal{K}[y, \theta]$  as average value of difference between models predictions. Therefore, this metric shows us how models in ensemble conform

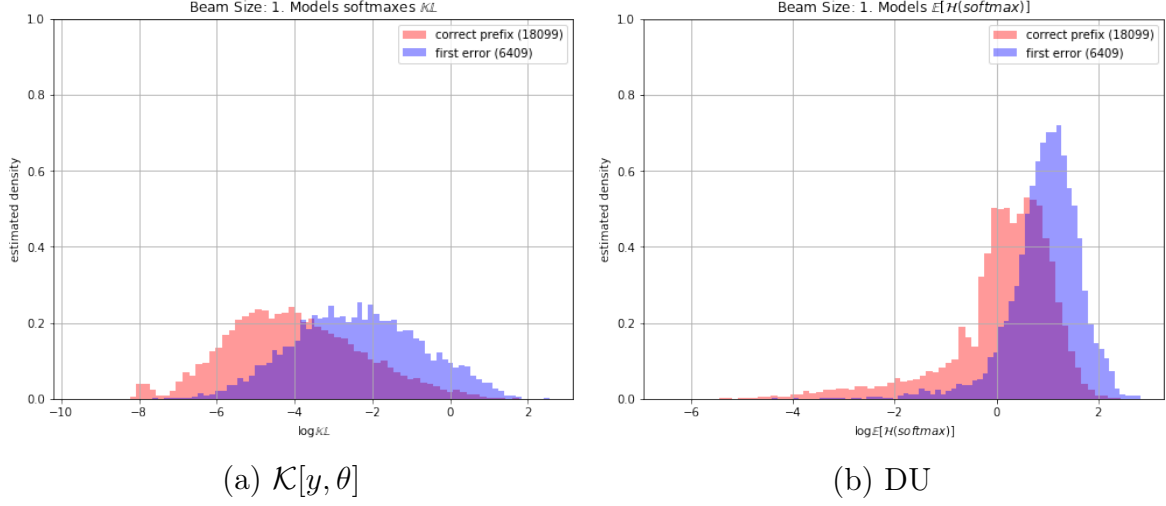


Figure 8 – Uncertainty estimates distribution for correct and error tokens

with each other. Prerequisite to use this kind of uncertainty estimate is possibly, the more inconsistent prediction models in ensemble make, the higher probability we get to make a mistake on current token. Less formal,  $\mathcal{K}[y, \theta]$  correlate with errors in predictions.

Final statistics can be seen at Fig. 8a. As in previous case, correct prefix density separable from first error token distribution. Moreover, expected value of kullback-leibler divergence is less than first error token's one. It means for us, that correlation we discussed earlier exists and it could be effective to use this kind of uncertainty estimate in misclassification problem.

Our final uncertainty estimate in this section is:

$$DU = \mathbb{E}_{q(\theta)}[\mathcal{H}[P(y|\mathbf{x}, \theta)]]$$

While we evaluate data uncertainty estimate, we consider average single model self-confidence in ensemble. This metric is an attempt to evaluate how common considered situation ( $p(y_i|y_{<i}, \mathbf{x}, \theta)$ ) is, how presented data blends with average model generalization ability. Therefore, DU could be a good measure of data uncertainty. Prerequisite to use DU is same as previous.

Final statistics for DU can be seen at Fig. 8b. We can see that our prerequisite is confirmed.

Overall, we considered three types of uncertainty estimation and empirically

proved correlation between all types of uncertainty and predictions errors. But anyway our main barrier to use uncertainty estimates in misclassification problem is token correctness indicator.

## 8.2 Using on practice

Let's assume that we were able to achieve impressive results by constructing some uncertainty estimate *PerfectTU*. How can we use this to improve quality?

All we can do is to mark positions in the sentences where we are wrong and, as a result, change our predictions. Due to the consistency of machine translation models, it would be logical to build a re-prediction not based on the original model, but to refer to some Oracle that produces more reliable predictions. In practice, such an Oracle could be some assessor or moderator. However, this approach imposes strong restrictions on the use of uncertainty estimation, since not every project can afford the support of human resources. Moreover, if you need to make fast online predictions in a task, you can no longer use moderators.

Thus, we need to change the way we evaluate uncertainty. We need to find a way of estimating so that our model can change the predictions based on the uncertainty estimates itself. Now this is not possible, because the predicting tokens with a fixed prefix have the same uncertainty rating. As our next approach, We do not build estimates for the entire softmax at once, but for each token in the softmax separately.

In the previous sections we have considered KL as an evaluation of knowledge uncertainty. Let us introduce a metric that correlates with KL, but is built for each token.

## 8.3 Beam search

We conduct missclassification detection experiment for different beam size (1, 5, 9, 20, 50). We explore that for both correctness indicator with beam size growth PR-AUC degradate. It is not a significant result, because the reason of

this phenomena lies in correctness indicator definition. For examle, if we consider  $correct(y^i, \mathbf{r})$ , while beam size grows, probability to meet incorrect token in beam search grows, too, because of low positive class ratio ( $\sim 0.25$ ).

Moreover, we do not find any correlation between beam size and uncertainty. As follow, considered uncertainty estimates can not solve larger beam size problem explicitly (in the way we do uncertainty estimation).

## 9 In-ensemble uncertainty estimate

Recall that for each token, we have denoted probability:

$$P(y_i|y_{<i}, x, \mathcal{D}) = \frac{1}{K} \sum_{j=1}^K P(y_i|y_{<i}, x, \theta_j)$$

Thus, for each token, we have: a conditional probability of the token for the ensemble and K conditional probabilities of the token for the models. Let us introduce an unbiased sample variance of the token probability along the models in the ensemble as a measure of uncertainty:

$$inens\_var(y_i) := S^2(P(y_i|y_{<i}, x, \theta_1), \dots, P(y_i|y_{<i}, x, \theta_K)) \quad (9)$$

This estimate serves as a measure of model consistency in predicting the probability of a given token (some analog of KL in previous chapters). Note that the entropy criterion is more suitable for us as a measure of uncertainty, since using variance it is difficult to separate cases: a singular distribution and a uniform distribution. However, variance still allows you to separate distributions with a high level of competing hypotheses. We choosed variance, because the probabilities of models do not form a probability measure (no distribution is set on these probabilities, the values may be completely inconsistent), so the calculation of entropy is incorrect in this case.

Consider the following graph (Fig. 9a). Here is a model of a new uncertainty

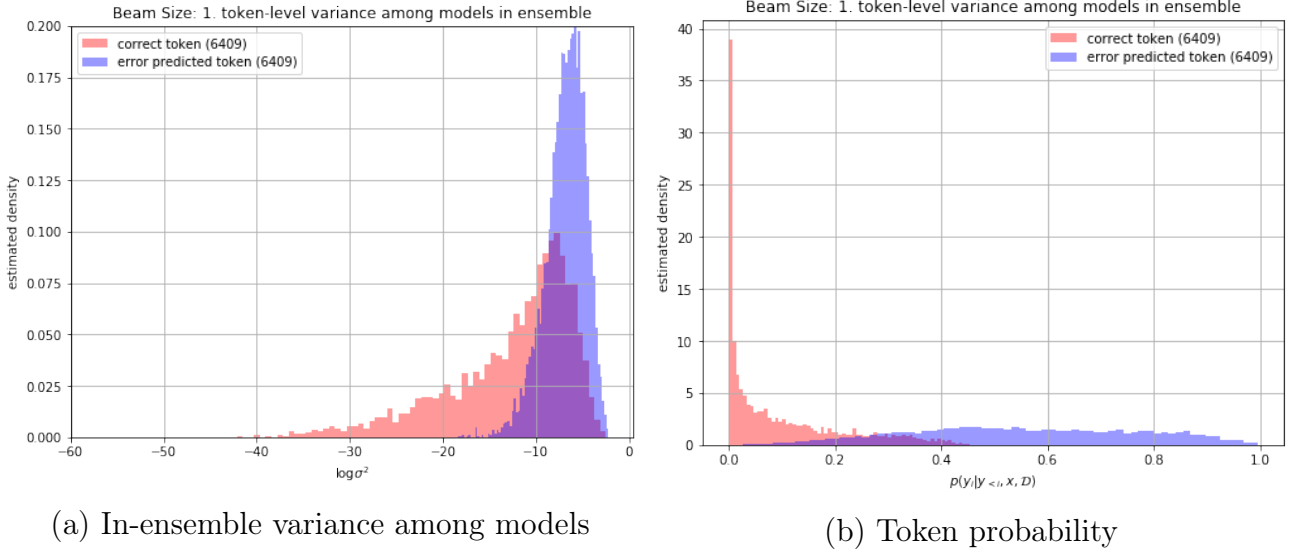


Figure 9 – Statistics for correct token and first error token in softmax

estimation for two types of tokens: the first false token in the prediction (the prefix before this token was predicted correctly) and the true token for this position with the same prefix. You can immediately notice that, as in the previous case, most of the probability mass of the correct token has lower variance than that of the incorrect token. This observation gives us hope that the models in the ensemble produce more consistent predictions of probabilities for the correct token and more doubt about the probability prediction for the final token.

There is a logical question: why, despite the lower consistency, the final prediction turned out to be the token that turned out to be. To answer this question, we need to look at the probability distributions for these same tokens (Fig. 9b).

As we can see from the charts, we make an error in this position, because often we have about zero probability for a true token, so we prefer the more likely token (and its continuation).

We would also like to see how the variance values and probability values relate to each other. Refer to Fig. 10. This graph reveals a serious disadvantage of considering variance. Consider the following example. Let us have some sample  $t = (t_1, \dots, t_m)$  and let it correspond to the sample variance:  $S^2(t)$ . Calculate:

$$S^2(\alpha t) = S^2((\alpha t_1, \dots, \alpha t_m)) = \alpha^2 S^2(t)$$

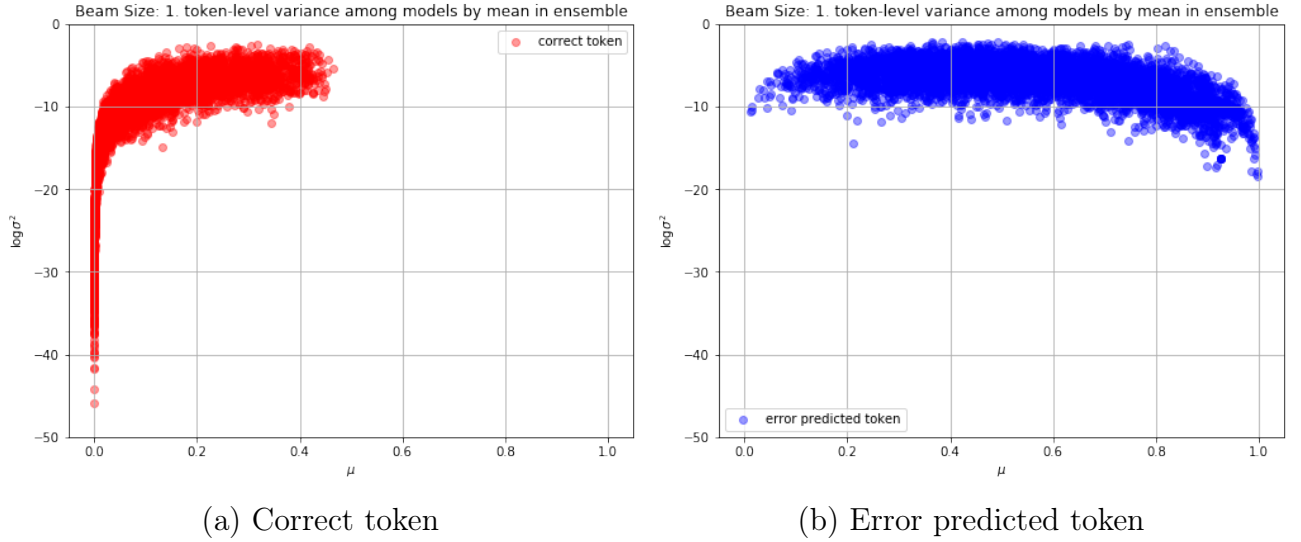


Figure 10 – Token-level ensemble probability by in-ensemble variance for correct token and first error token

Thus, the value of our metric depends not only on the total relation between values in the sample, but also on the absolute values of these values. This is the effect we saw on chart 10. For a higher probability mass of a valid token, the low value of the uncertainty estimate that we selected is not due to the consistency of the models, but rather to the values of probabilities, i.e. the variance of near zero probabilities always lower, because of property above.

However, despite this, by comparing the variance of the correct and false token with equal probabilities, we can observe that the expected variance for the correct token is lower than for the false one. Therefore, although not on the entire probability mass, but on some of its part, our estimation of uncertainty gives adequate predictions.

In order to make sure of this, you need to see how our rating behaves on other tokens of the same softmax. To do this, we consider: a random token, tokens with top-3, top-5, and top-20 probabilities in a descending series of variations. The final graphs of the probability distribution and uncertainty estimates can be seen in Fig. 11.

As we can see, correct tokens are not so hopeless. A high proportion of the probability mass of correct tokens has relatively high probabilities (the expected position in the variation series of probabilities == 2). Moreover, the models ’

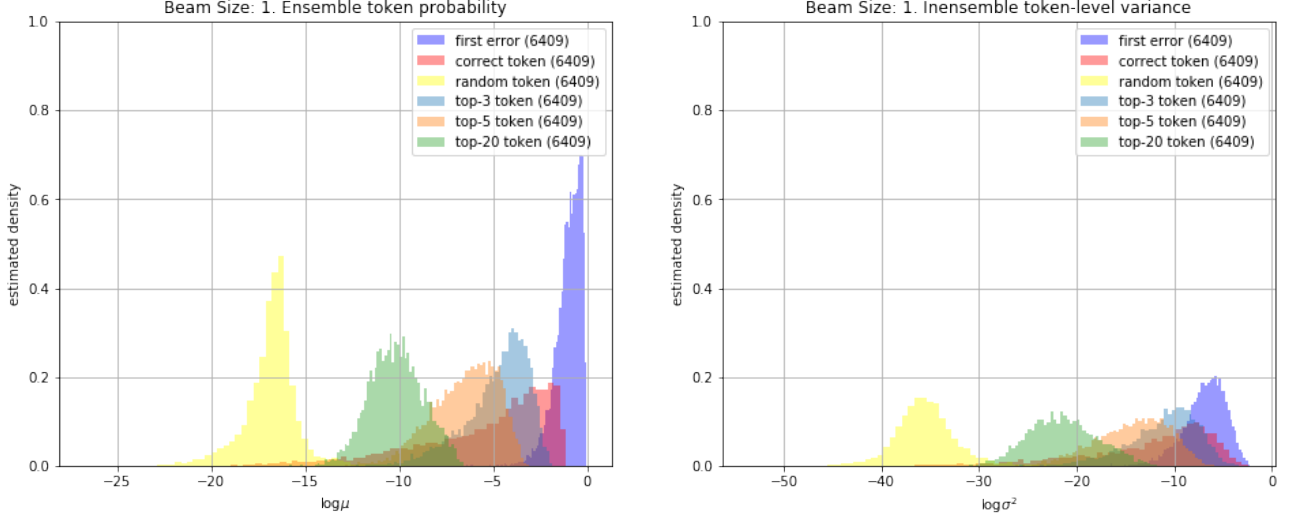


Figure 11 – Token-level ensemble probability and in-ensemble variance for correct token, first error token, random token, top-3 top-5 and top-20 tokens

confidence in the correct tokens is not inferior to the variance of other tokens. We can use these observations to help the model find the correct tokens.

## 9.1 Probability distribution calibration

Our task now is to renormalize the probability distributions of tokens with the uncertainty metric in such a way that the distribution of the correct token becomes more prevalent in terms of choosing the final hypothesis by the model. Thus, the problem is reduced to finding some function  $g$ :

$$\hat{p}(y|x, \mathcal{D}) = \text{softmax}(g(p(y|x, \mathcal{D}), \text{inens\_var}(y)))$$

If we look at Fig.. 12, it becomes clear that by selecting  $g$ , we can only bend this graph along the x axis. In visual terms, we would like to pull the density of the correct token to the right edge along each contour line of the variance, and the density of any other token to the left.

Remind that we want to penalize for high values of uncertainty, but not so much as to overweight completely irrelevant tokens. A bad example of the  $g$  function would be:



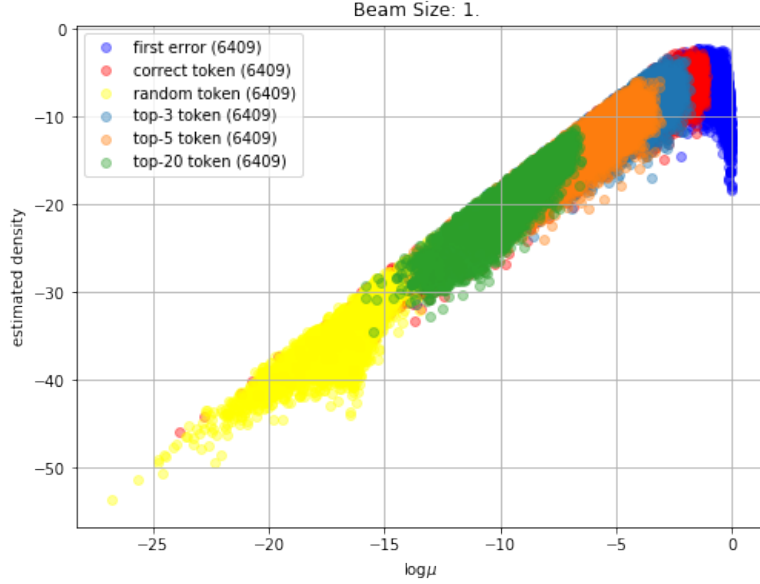


Figure 12 – Token-level ensemble probability by in-ensemble variance for correct token, first error token, random token, top-3 top-5 and top-20 tokens

Here and further, let us denote:  $\sigma^2 := inens\_var$

$$g_{naive} := \log p(y|x, \mathcal{D}) - \log(\sigma^2) \quad (10)$$

We explicitly reward tokens with a very low variance value and, conversely, heavily penalize them for high uncertainty values. This naive approach leads to a diametrically opposite situation that we want to get (Fig. 13a). We can see that such strong penalties lead to a high probability of accepting random tokens. This is because the variance of random tokens is extremely low due to the property of the variance that we discussed earlier. As a result, the model simply produces completely random translations.

Let us try to slightly weaken the penalty for uncertainty, divide the probability by the root of the variance:

$$g_{less\_naive} := \log p(y|x, \mathcal{D}) - \log(\sigma) \quad (11)$$

As we can see from Fig. 13b, this did not bring any success, since we did not even change the relative order of probability distributions.

If we take another look at Fig. 12, it is clear that we cannot use linear

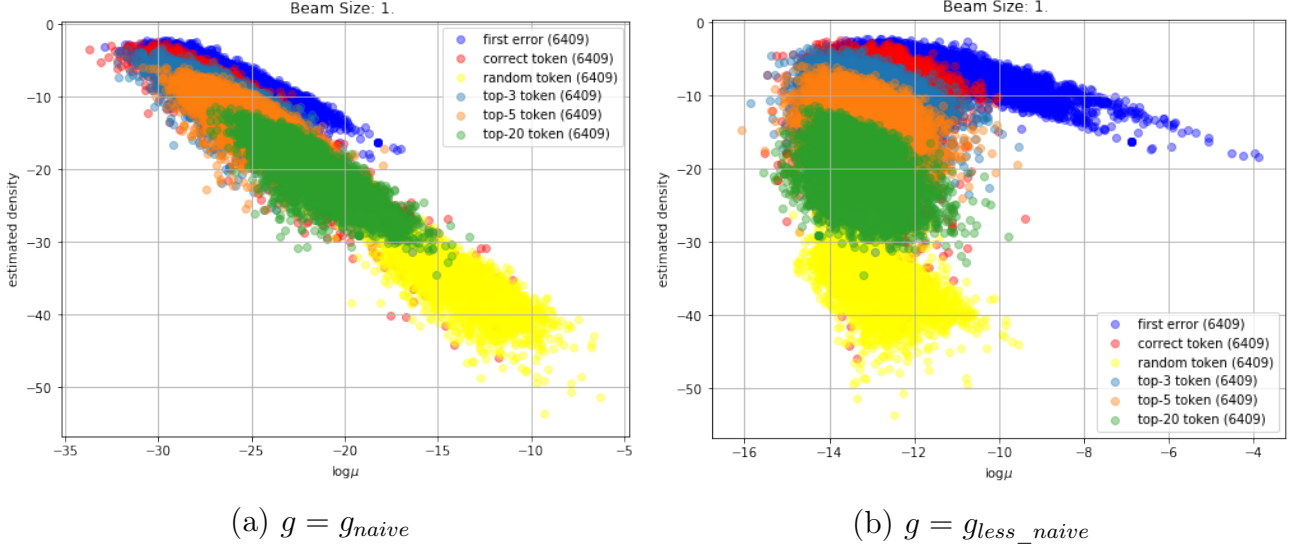


Figure 13 – Token-level ensemble probability by in-ensemble variance for correct token, first error token, random token, top-3 top-5 and top-20 tokens for  $g$

functions  $g$  to rearrange the order of distributions in terms of probabilities (recall that the graph is presented in log scales). Thus, we come to the conclusion that the calibration must not be linear by  $\log \sigma^2$ , moreover, the peak (maximum point) must fall on a part of the distribution of the correct token, which is not overlapped by any other distribution.

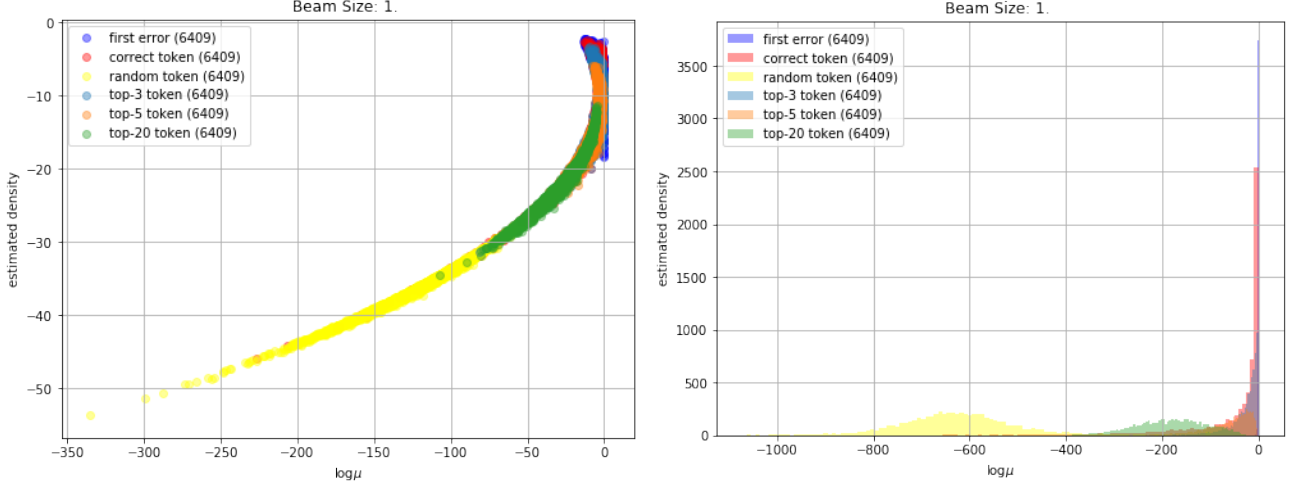
Then let us take a function:

$$g_2 := \log p(y|x, \mathcal{D}) - (\log(\sigma) + \alpha)^2 \quad (12)$$

Where  $\alpha$  is a certain hyperparameter that depends on which point the maximum is reached. Let us consider a special example (Fig. 14a).

As we can see, now most probability mass of correct token moved forward to first positions in probability variational series, but also top tokens moved forward too. By the reason of fact, all tokens' distributions are mixed, we can not certainly separate correct token distribution from any other. But we can make correct token more prevailing over any other token on a fixed variance counter level. As we can see on a chart, we achieve this goal.

Note that the density of the wrong token still remains to the right of the density of the correct token. Unfortunately or fortunately, we can't do anything



(a) Calibrated probabilities by variance

(b) Calibrated probabilities histogram

Figure 14 –  $g_2$  statistics for correct token, first error token, random token, top-3 top-5 and top-20 tokens.

$$g_2 = \log p(y|x, \mathcal{D}) - (\log(\sigma) + 4)^2$$

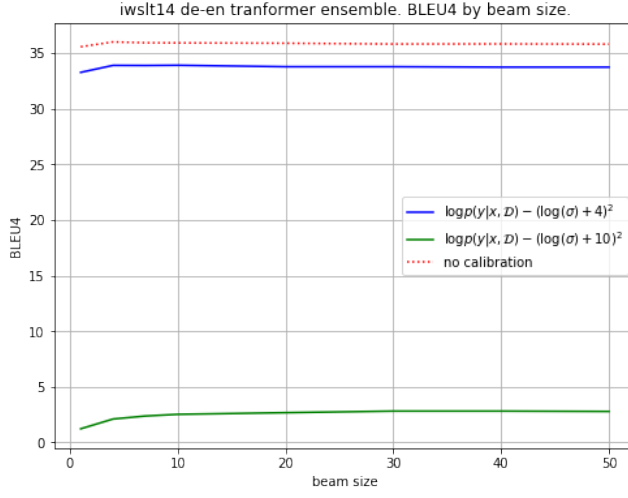
with these points, because if a certain token has a high probability on the same contour line of variance, obviously, we should give preference to this particular token. However, even this results can lead to good results, because in the case of a non-single beam size, even a small overweight can lead to the choice of correct hypotheses, because the choice of hypothesis is based on cumulative probabilities, not on a greedy choice.

If we take a look at Fig. 14b, we can see ensemble probability distributions after  $g_2$  calibration. Our hypothesis about correct token prevailation confirms.

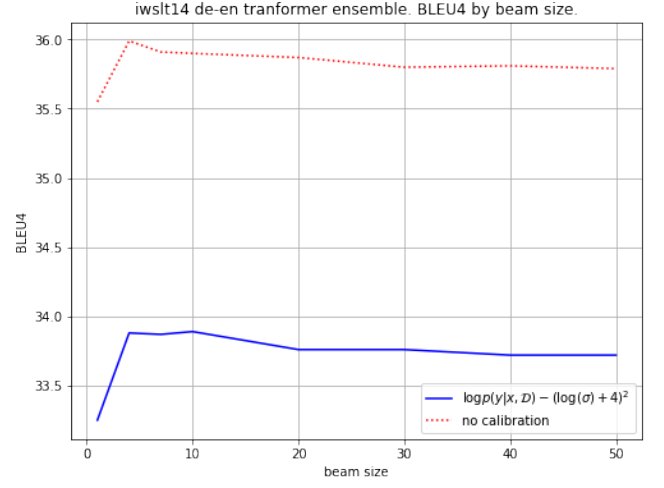
An important note about results is that, while we use beam search with beam size equal 1, we do not earn better results using  $g_2$  calibration, because most part of first error token is still most probable. But if we consider beam search model with beam size more than 1, we can expect quality improvement, because non-greedy beam search not always chooses most likely token and estimate next token probability using cumulative prefix probability.

## 9.2 Calibration BLEU

Let us consider BLEU4 evaluation on a test set for calibrated ensemble (Fig. 15).

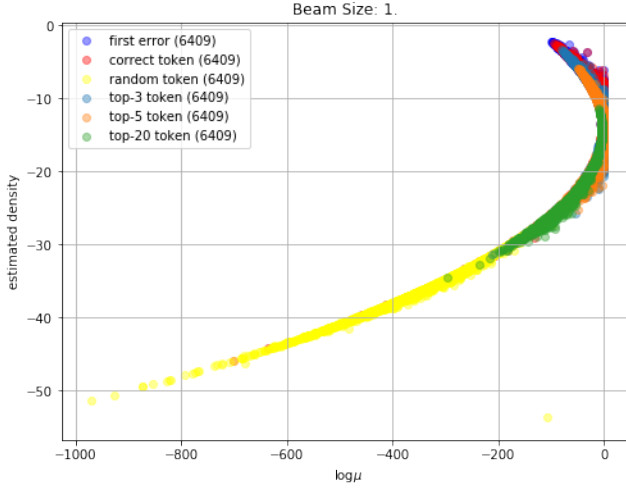


(a)

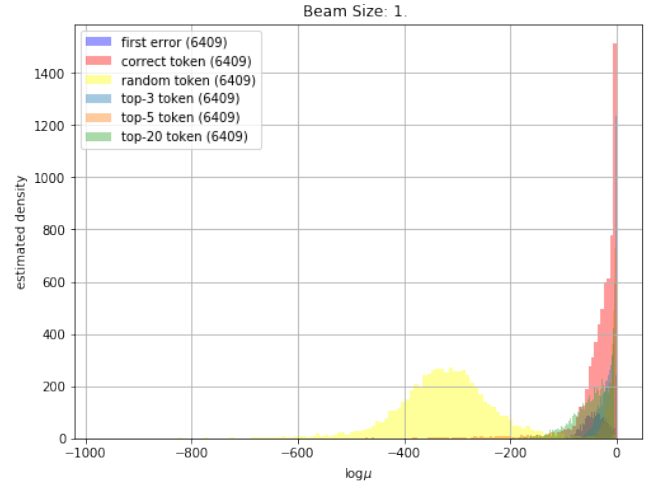


(b)

Figure 15 – BLEU4 by beam size for calibrated transformer ensemble



(a) Calibrated probabilities by variance



(b) Calibrated probabilities histogram

Figure 16 –  $g_2$  statistics for correct token, first error token, random token, top-3 top-5 and top-20 tokens.

$$g_2 = \log p(y|x, \mathcal{D}) - (\log(\sigma) + 10)^2$$

As we can see, baseline ensemble without any calibration overperforms all our tested calibrated models. The worst results are presented by square calibration with bias equal 10. BLEU is less than 5 for this case. The reason of this poor results can be seen at Fig. 16. After calibration all top-20 (including all more probable tokens) mix together that is why in lots of cases softmaxes have more entropy (recall it means that our model more uncertain in predictions), so we get a lot of more random predicted tokens. The problem is we choose bias too large and calibration scatter peak is not on correct token distribution but on irrelevant tokens (top-5 top-20 tokens). Therefore, we get too low BLEU.

More significant results are presented by calibration with bias equal 4. The problem in this case is similar to previous one. Now peak is placed on more relevant tokens, but still irrelevant (top-3) tokens enter in resulting predictions. Therefore, let us conduct same experiment but with less bias.

Overall, we consider different kinds of calibration, but we did not achieve quality improvement. The reason of poor final results lie on a weak uncertainty estimate. If we look at Fig. 11, we can see that in-ensemble variance distribution for correct token overlap all tokens distributions (except random one) that is why, it is too hard to separate correct token from others. Too small part of correct token distribution for variance is separable from others. This part is too small to get enough quality improvement.

## 10 Conclusion

In this work we conduct misclassification detection experiments, using different types of uncertainty estimates, on neural machine translation task. We explore correlation between uncertainty and predictions error, analysing uncertainty estimates distributions. Our conclusion is that it is possible to detect prediction error via uncertainty estimation, because considered uncertainty estimates correlate with errors (correct tokens have lower expected values than others). At the same time, it is hard to evaluate correct quality estimation for misclassification detection via uncertainty, because there are lots of problems with defining reasonable tokens correctness measure.

Moreover, we go deeper softmax level and explore new way of uncertainty estimation. We introduce in-ensemble uncertainty estimate for every token. We approve correlation between new measure and predictions error. As well, we try to recalibrate token probabilities using token-level in-ensemble variance.

Unfortunately, we do not achieve model quality improvement, but also we do not claim that this hypothesis is wrong. Further studies should conduct same experiments using different uncertainty estimates. As we conclude in-ensemble

variance distribution has weak correlation with error predictions to achieve high quality improvement.

## References

1. Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, Yann N. Dauphin. Facebook AI Research. Convolutional Sequence to Sequence Learning. PMLR 70, 2017.
2. Ilya Sutskever, Oriol Vinyals, Quoc V. Le. Sequence to Sequence Learning with Neural Networks. Google. NIPS 2014.
3. Dzmitry Bahdanau, KyungHyun Cho, Yoshua Bengio. Neural Machine Translation by Jointly Learning to Align and Translate. ICLR 2015.
4. Kenton Murray, David Chiang. Department of Computer Science and Engineering, University of Notre Dame. Correcting Length Bias in Neural Machine Translation. WMT 2018.
5. Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi. Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation. ArXiv 2016
6. Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser. Attention Is All You Need. NIPS 2017.
7. Philipp Koehn, Rebecca Knowles. Six Challenges for Neural Machine Translation. NMT@ACL 2017.
8. Aviral Kumar, Sunita Sarawagi. Calibration of Encoder Decoder Models for Neural Machine Translation. ArXiv 2019.
9. Myle Ott, Michael Auli, David Grangier, Marc'Aurelio Ranzato. Analyzing Uncertainty in Neural Machine Translation. PMLR 80, 2018.

10. Andrey Malinin, Mark Gales. Predictive Uncertainty Estimation via Prior Networks. NeurIPS 2018.
11. Andrey Malinin, Mark Gales. Uncertainty in Structured Prediction. 2020.
12. Facebook. AI Research Sequence-to-Sequence Toolkit
13. IWSLT 2014, The International Workshop on Spoken Language Translation