

# Language Model Prior for Low-Resource Neural Machine Translation

Christos Baziotis, Barry Haddow, Alexandra Birch

Institute for Language, Cognition and Computation  
School of Informatics, University of Edinburgh  
Edinburgh, UK

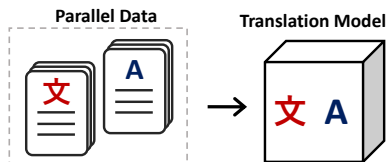
EMNLP 2020



## Neural Machine Translation

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y}} \log p(\mathbf{y}|\mathbf{x})$$

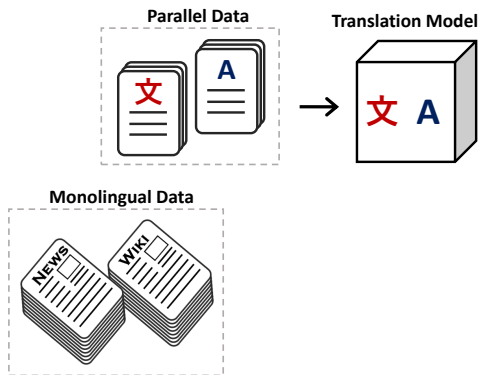
Modeling directly  $p(\mathbf{y}|\mathbf{x})$  requires **large** amounts of parallel data.



## Neural Machine Translation

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y}} \log p(\mathbf{y}|\mathbf{x})$$

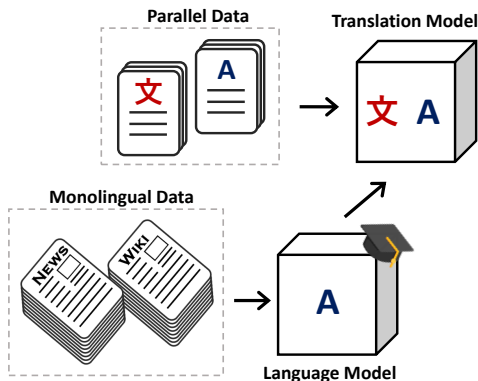
Modeling directly  $p(\mathbf{y}|\mathbf{x})$  requires **large** amounts of parallel data.



## Neural Machine Translation

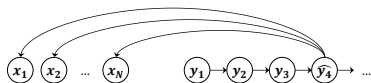
$$\hat{y} = \arg \max_y \log p(y|x)$$

Modeling directly  $p(y|x)$  requires **large** amounts of parallel data.



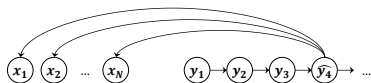
## Noisy Channel

It models the “reverse translation probability”  $p(\mathbf{y}|\mathbf{x}) \propto p_{\text{TM}}(\mathbf{x}|\mathbf{y}) \times p_{\text{LM}}(\mathbf{y})$



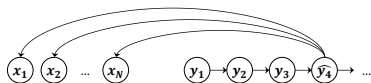
## Noisy Channel

It models the “reverse translation probability”  $p(\mathbf{y}|\mathbf{x}) \propto p_{\text{TM}}(\mathbf{x}|\mathbf{y}) \times p_{\text{LM}}(\mathbf{y})$



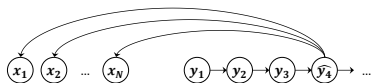
## Noisy Channel

It models the “reverse translation probability”  $p(\mathbf{y}|\mathbf{x}) \propto p_{\text{TM}}(\mathbf{x}|\mathbf{y}) \times p_{\text{LM}}(\mathbf{y})$



## Noisy Channel

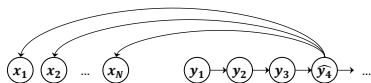
It models the “reverse translation probability”  $p(\mathbf{y}|\mathbf{x}) \propto p_{\text{TM}}(\mathbf{x}|\mathbf{y}) \times p_{\text{LM}}(\mathbf{y})$



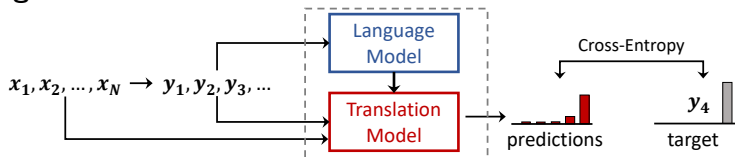


## Noisy Channel

It models the “reverse translation probability”  $p(\mathbf{y}|\mathbf{x}) \propto p_{\text{TM}}(\mathbf{x}|\mathbf{y}) \times p_{\text{LM}}(\mathbf{y})$

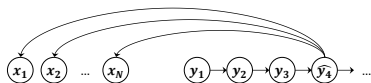


## Language Model Fusion

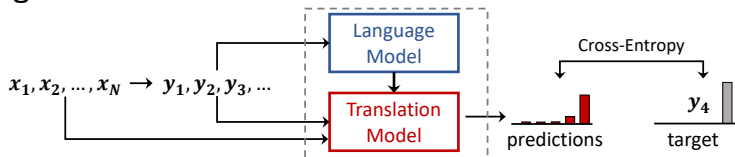


## Noisy Channel

It models the “reverse translation probability”  $p(\mathbf{y}|\mathbf{x}) \propto p_{\text{TM}}(\mathbf{x}|\mathbf{y}) \times p_{\text{LM}}(\mathbf{y})$



## Language Model Fusion

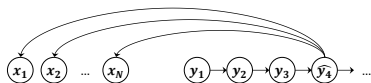


**Feature-based:** incorporate the *representations* of the LM

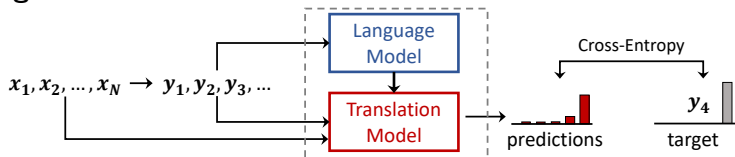
- deep-fusion (Gulcehre et al., 2015), cold-fusion (Sriram et al., 2018)

## Noisy Channel

It models the “reverse translation probability”  $p(\mathbf{y}|\mathbf{x}) \propto p_{\text{TM}}(\mathbf{x}|\mathbf{y}) \times p_{\text{LM}}(\mathbf{y})$



## Language Model Fusion



**Feature-based:** incorporate the *representations* of the LM

- deep-fusion (Gulcehre et al., 2015), cold-fusion (Sriram et al., 2018)

**Probability Interpolation:** combine the *probabilities* of the TM & the LM

- shallow-fusion (Gulcehre et al., 2015), simple-fusion (Stahlberg et al., 2018)

# Limitations of Prior Work

- 1 **Computational overhead**, as the LM is also required for decoding

# Limitations of Prior Work

- 1 **Computational overhead**, as the LM is also required for decoding
- 2 **Unwanted LM-TM interference** in probability interpolation

# Limitations of Prior Work

- 1 **Computational overhead**, as the LM is also required for decoding
- 2 **Unwanted LM-TM interference** in probability interpolation

$$p(\mathbf{y}_t) = \text{softmax}(\log p_{\text{TM}}(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{x}) + \log p_{\text{LM}}(\mathbf{y}_t | \mathbf{y}_{<t})) \quad \text{POSTNORM}$$

(simple-fusion; Stahlberg et al., 2018)

**DE:** die Republikaner im Kongress drängen auf eine umfassendere Neufassung der Ozonregeln.

**EN:** Republicans → in → Congress → are → pushing → for → a → broader → rewrite → of → the → ozone → rules.

TM		LM	
broader	34.3%	new	8.9%
wider	22.4%	repeal	7.1%
larger	7.2%	bill	3.7%
...		...	

× =

more	44.8%
wider	31.8%
larger	11.8%
...	

Real example of POSTNORM from DE→EN newstest2018

# Limitations of Prior Work

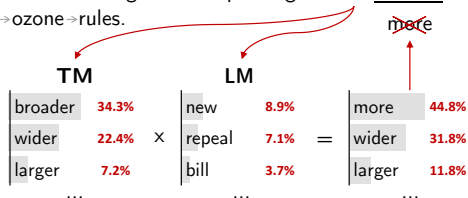
- 1 **Computational overhead**, as the LM is also required for decoding
- 2 **Unwanted LM-TM interference** in probability interpolation

$$p(\mathbf{y}_t) = \text{softmax}(\log p_{\text{TM}}(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{x}) + \log p_{\text{LM}}(\mathbf{y}_t | \mathbf{y}_{<t})) \quad \text{POSTNORM}$$

(simple-fusion; Stahlberg et al., 2018)

**DE:** die Republikaner im Kongress drängen auf eine umfassendere Neufassung der Ozonregeln.

**EN:** Republicans → in → Congress → are → pushing → for → a → broader → rewrite → of → the → ozone → rules.



Real example of POSTNORM from DE → EN newstest2018

# Limitations of Prior Work

- 1 **Computational overhead**, as the LM is also required for decoding
- 2 **Unwanted LM-TM interference** in probability interpolation

$$p(\mathbf{y}_t) = \text{softmax}(\log p_{\text{TM}}(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{x}) + \log p_{\text{LM}}(\mathbf{y}_t | \mathbf{y}_{<t})) \quad \text{POSTNORM}$$

(simple-fusion; Stahlberg et al., 2018)

**DE:** die Republikaner im Kongress drängen auf eine umfassendere Neufassung der Ozonregeln.

**EN:** Republicans → in → Congress → are → pushing → for → a → broader → rewrite → of → the → ozone → rules.

TM		LM	
broader	34.3%	new	8.9%
wider	22.4%	repeal	7.1%
larger	7.2%	bill	3.7%
...		...	

× =

more	44.8%
wider	31.8%
larger	11.8%
...	

Real example of POSTNORM from DE→EN newstest2018



# Limitations of Prior Work

- 1 **Computational overhead**, as the LM is also required for decoding
- 2 **Unwanted LM-TM interference** in probability interpolation

$$p(\mathbf{y}_t) = \text{softmax}(\log p_{\text{TM}}(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{x}) + \log p_{\text{LM}}(\mathbf{y}_t | \mathbf{y}_{<t})) \quad \text{POSTNORM}$$

(simple-fusion; Stahlberg et al., 2018)

**DE:** die Republikaner im Kongress drängen auf eine umfassendere Neufassung der Ozonregeln.

**EN:** Republicans → in → Congress → are → pushing → for → a → broader → rewrite → of → the → ozone → rules.

TM		LM	
broader	34.3%	new	8.9%
wider	22.4%	repeal	7.1%
larger	7.2%	bill	3.7%
...		...	

× =

more	44.8%
wider	31.8%
larger	11.8%
...	

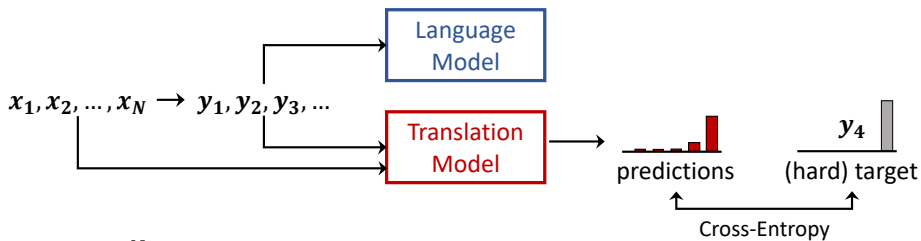
Real example of POSTNORM from DE → EN newstest2018

# Language Model Prior

We regularize the TM, by using the LM as **prior** over the decoder

# Language Model Prior

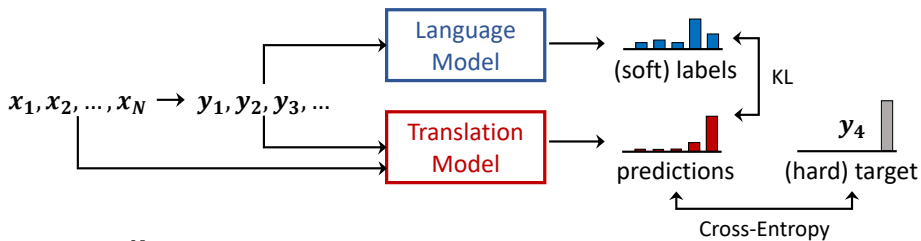
We regularize the TM, by using the LM as **prior** over the decoder



$$\mathcal{L} = \sum_{t=1}^N \underbrace{-\log p_{\text{TM}}(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{x})}_{\text{translation}}$$

# Language Model Prior

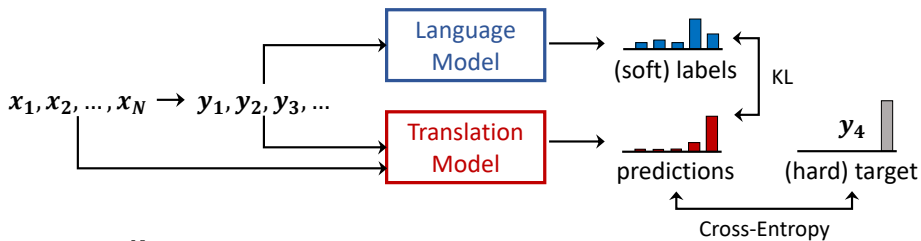
We regularize the TM, by using the LM as **prior** over the decoder



$$\mathcal{L} = \sum_{t=1}^N \underbrace{-\log p_{\text{TM}}(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{x})}_{\text{translation}} + \lambda \underbrace{D_{\text{KL}}(p_{\text{LM}}(\mathbf{y}_t | \mathbf{y}_{<t}) || p_{\text{TM}}(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{x}))}_{\text{prior}}$$

# Language Model Prior

We regularize the TM, by using the LM as **prior** over the decoder



$$\mathcal{L} = \sum_{t=1}^N \underbrace{-\log p_{\text{TM}}(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{x})}_{\text{translation}} + \lambda \underbrace{D_{\text{KL}}(p_{\text{LM}}(\mathbf{y}_t | \mathbf{y}_{<t}) \parallel p_{\text{TM}}(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{x}))}_{\text{prior}}$$

- + The TM can **overrule** the LM when needed
- + **Easier** than priors on neural network weights
- + The LM is **not required** during decoding

# Relation to Other Methods with Soft Targets

## Knowledge Distillation

Also uses **soft targets** from a teacher model, but needs **parallel data**

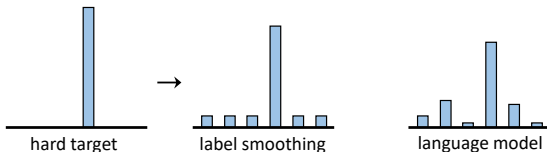
# Relation to Other Methods with Soft Targets

## Knowledge Distillation

Also uses **soft targets** from a teacher model, but needs **parallel data**

## Label Smoothing

Converts the one-hot into **soft targets** to penalize **over-confidence**



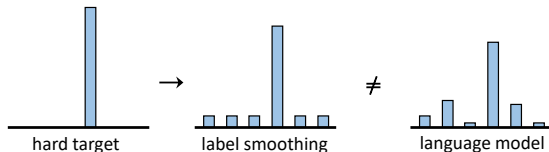
# Relation to Other Methods with Soft Targets

## Knowledge Distillation

Also uses **soft targets** from a teacher model, but needs **parallel data**

## Label Smoothing

Converts the one-hot into **soft targets** to penalize **over-confidence**



- Assigns **equal** probability to all the incorrect classes, unlike the LM.



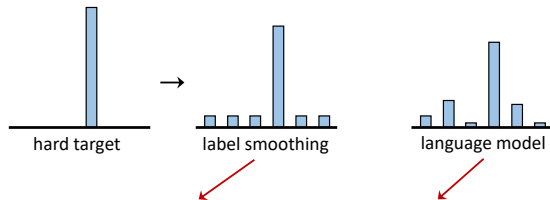
# Relation to Other Methods with Soft Targets

## Knowledge Distillation

Also uses **soft targets** from a teacher model, but needs **parallel data**

## Label Smoothing

Converts the one-hot into **soft targets** to penalize **over-confidence**



$$\mathcal{L} = \sum_{t=1}^N \underbrace{-\log p_{\text{TM}}(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{x})}_{\text{translation}} + \lambda \underbrace{D_{\text{KL}}(p_{\text{LM}}(\mathbf{y}_t | \mathbf{y}_{<t}) \parallel p_{\text{TM}}(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{x}))}_{\text{prior}}$$

- Assigns **equal** probability to all the incorrect classes, unlike the LM.
- **Orthogonal** to posterior regularization (LM-prior)

# Experimental Setup

Language Pair	Sentences	Corpus
English-German	275K	News Commentary v13
English-Turkish	190K	SETIMES

Table: Parallel data from WMT-2018

Language	Small (3M)	Large (30M)	Corpus
English	✓	✓	News Crawl 2016
German	✓	✓	News Crawl 2016
Turkish	✓	–	News Crawl 2010-2018

Table: Monolingual data from News Crawl articles

## Models

- Architecture: Transformer (LMS and TMS)
- Vocabulary: 16K symbols (sentencepiece)

# Translation Results

Method	DE→EN	EN→DE	TR→EN	EN→TR
Base	26.6	25.6	16.6	11.2
Shallow-fusion	27.8	26.0	17.3	11.5
POSTNORM+ LS	26.4	23.3	16.0	11.0
Base + LS	28.4	27.3	18.4	12.6
Base + Prior	<b>30.2</b>	<b>29.1</b>	<b>19.5</b>	<b>13.8</b>

**Table:** BLEU scores of each NMT method (mean 3 runs)

# Translation Results

Method	DE→EN	EN→DE	TR→EN	EN→TR
Base	26.6	25.6	16.6	11.2
Shallow-fusion	27.8	26.0	17.3	11.5
POSTNORM+ LS	26.4	23.3	16.0	11.0
Base + LS	28.4	27.3	18.4	12.6
Base + Prior	<b>30.2</b>	<b>29.1</b>	<b>19.5</b>	<b>13.8</b>

Table: BLEU scores of each NMT method (mean 3 runs)

# Translation Results

Method	DE→EN	EN→DE	TR→EN	EN→TR
Base	26.6	25.6	16.6	11.2
Shallow-fusion	27.8	26.0	17.3	11.5
POSTNORM+ LS	26.4	23.3	16.0	11.0
Base + LS	28.4	27.3	18.4	12.6
Base + Prior	<b>30.2</b>	<b>29.1</b>	<b>19.5</b>	<b>13.8</b>

Table: BLEU scores of each NMT method (mean 3 runs)

# Translation Results

Method	DE→EN	EN→DE	TR→EN	EN→TR
Base	26.6	25.6	16.6	11.2
Shallow-fusion	27.8	26.0	17.3	11.5
POSTNORM+ LS	26.4	23.3	16.0	11.0
Base + LS	28.4	27.3	18.4	12.6
Base + Prior	<b>30.2</b>	<b>29.1</b>	<b>19.5</b>	<b>13.8</b>

Table: BLEU scores of each NMT method (mean 3 runs)

# Translation Results

Method	DE→EN	EN→DE	TR→EN	EN→TR
Base	26.6	25.6	16.6	11.2
Shallow-fusion	27.8	26.0	17.3	11.5
POSTNORM+ LS	26.4	23.3	16.0	11.0
Base + LS	28.4	27.3	18.4	12.6
Base + Prior	<b>30.2 (+1.8)</b>	<b>29.1 (+1.8)</b>	<b>19.5 (+1.1)</b>	<b>13.8 (+1.2)</b>

Table: BLEU scores of each NMT method (mean 3 runs)

# Translation Results

Method	DE→EN	EN→DE	TR→EN	EN→TR
Base	26.6	25.6	16.6	11.2
Shallow-fusion	27.8	26.0	17.3	11.5
POSTNORM+ LS	26.4	23.3	16.0	11.0
Base + LS	28.4	27.3	18.4	12.6
Base + Prior	<b>30.2</b>	<b>29.1</b>	<b>19.5</b>	<b>13.8</b>
Base + Prior + LS	<u>30.3</u> (+0.1)	<u>29.7</u> (+0.6)	19.5	<u>14.1</u> (+0.3)

Table: BLEU scores of each NMT method (mean 3 runs)

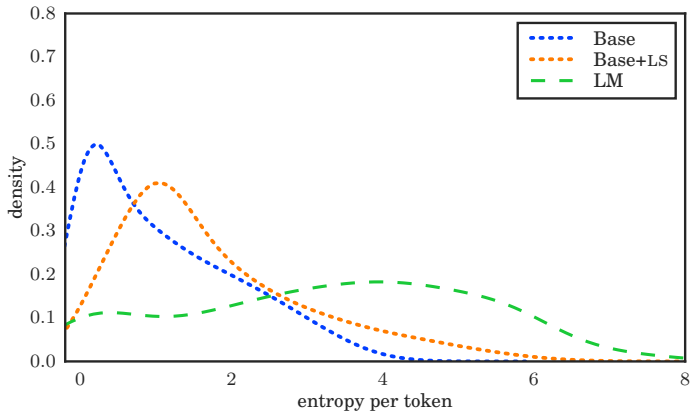


# Translation Results

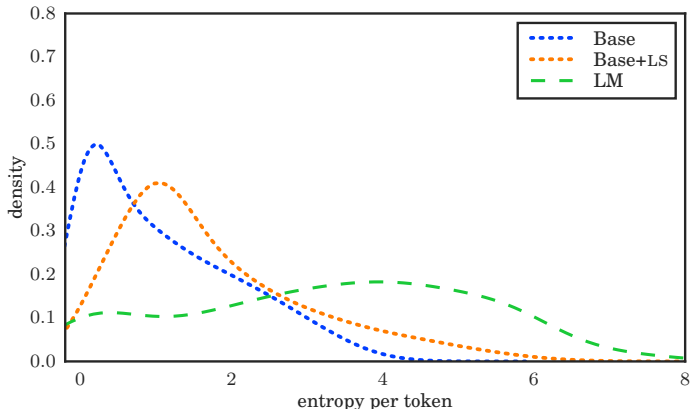
Method	DE→EN	EN→DE	TR→EN	EN→TR
Base	26.6	25.6	16.6	11.2
Shallow-fusion	27.8	26.0	17.3	11.5
POSTNORM+ LS	26.4	23.3	16.0	11.0
Base + LS	28.4	27.3	18.4	12.6
Base + Prior	<b>30.2</b>	<b>29.1</b>	<b>19.5</b>	<b>13.8</b>
Base + Prior + LS	<u>30.3</u>	<u>29.7</u>	19.5	<u>14.1</u>
Base + Prior (30M)	30.0	<u>29.8</u> (+0.7)	19.5	–

Table: BLEU scores of each NMT method (mean 3 runs)

# Analysis

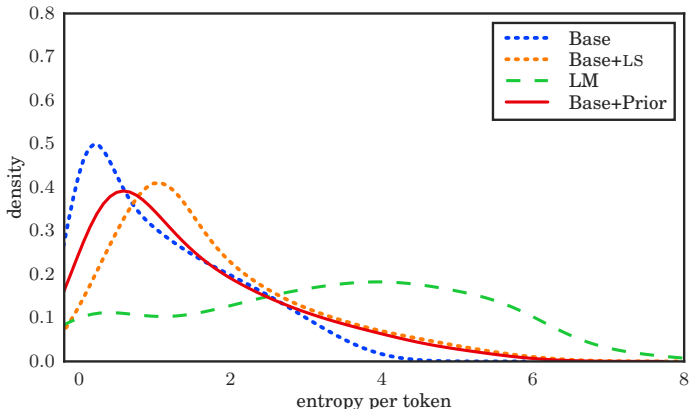


Estimated density of each model's entropy on the  $DE \rightarrow EN$  test set



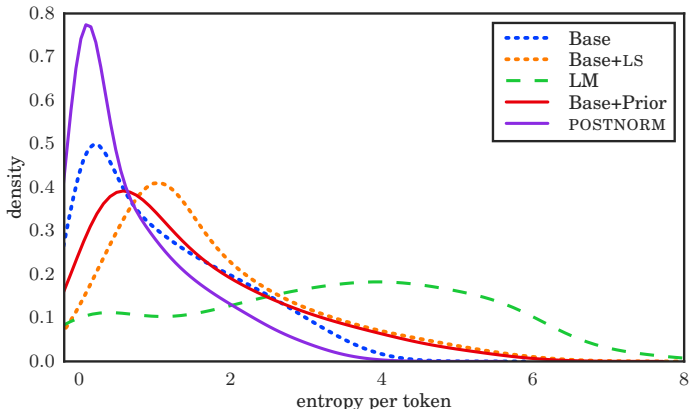
Estimated density of each model's entropy on the  $DE \rightarrow EN$  test set

- LS successfully mitigates overfitting, by **penalizing confidence**



Estimated density of each model's entropy on the  $DE \rightarrow EN$  test set

- LS successfully mitigates overfitting, by **penalizing confidence**
- “Base+Prior” does **not just** penalize confidence, but **exploits** the LM



Estimated density of each model's entropy on the  $DE \rightarrow EN$  test set

- LS successfully mitigates overfitting, by **penalizing confidence**
- “Base+Prior” does **not just** penalize confidence, but **exploits** the LM
- “POSTNORM” is **over-confident**, due to unwanted interference from LM

- 1 **Simple** approach to incorporate prior knowledge in NMT
  - The decoupling of the LM from the TM enables **fast decoding**
  - We change **only** the training objective
- 2 **Promising results** in two low-resource translation datasets
  - Improvements even with **modest** monolingual data
- 3 **Analysis** of TM output distributions using different methods
  - Evidence that the model **exploits** the **knowledge** of the LM-prior

Thank you

## Bonus Slides



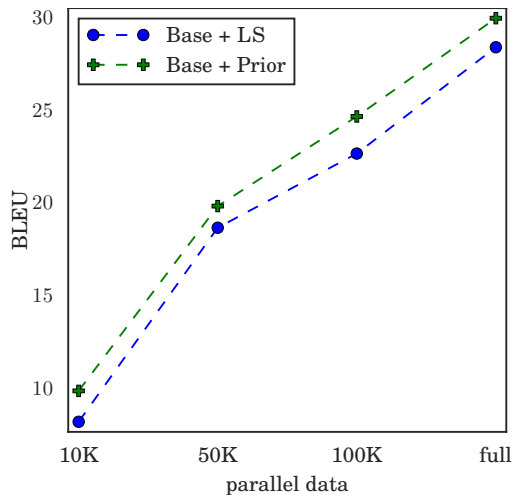
# Full Results

Method	DE→EN		EN→DE		TR→EN		EN→TR	
	dev	test	dev	test	dev	test	dev	test
Base	22.6 $\pm$ 0.1	26.6 $\pm$ 0.1	18.3 $\pm$ 0.3	25.6 $\pm$ 0.2	15.9 $\pm$ 0.0	16.6 $\pm$ 0.3	12.2 $\pm$ 0.1	11.2 $\pm$ 0.2
Shallow-fusion	23.4 $\pm$ 0.1	27.8 $\pm$ 0.1	18.5 $\pm$ 0.2	26.0 $\pm$ 0.1	16.5 $\pm$ 0.1	17.3 $\pm$ 0.3	12.7 $\pm$ 0.0	11.5 $\pm$ 0.1
POSTNORM	20.4 $\pm$ 0.2	24.5 $\pm$ 0.3	16.6 $\pm$ 0.1	22.9 $\pm$ 0.3	13.8 $\pm$ 0.2	11.0 $\pm$ 0.1	10.0 $\pm$ 0.1	10.2 $\pm$ 0.1
Base + LS	23.8 $\pm$ 0.6	28.4 $\pm$ 0.7	19.2 $\pm$ 0.3	27.3 $\pm$ 0.3	17.5 $\pm$ 0.1	18.4 $\pm$ 0.2	13.8 $\pm$ 0.2	12.6 $\pm$ 0.0
Base + Prior	<b>24.9<math>\pm</math>0.0</b>	<b>30.2<math>\pm</math>0.1</b>	<b>20.5<math>\pm</math>0.3</b>	<b>29.1<math>\pm</math>0.7</b>	<b>18.5<math>\pm</math>0.2</b>	<b>19.5<math>\pm</math>0.2</b>	<b>15.1<math>\pm</math>0.1</b>	<b>13.8<math>\pm</math>0.1</b>
Base + Prior + LS	<u>25.1<math>\pm</math>0.3</u>	<u>30.3<math>\pm</math>0.3</u>	<u>20.8<math>\pm</math>0.4</u>	<u>29.7<math>\pm</math>0.7</u>	18.5 $\pm$ 0.3	19.5 $\pm$ 0.2	<u>15.5<math>\pm</math>0.1</u>	<u>14.1<math>\pm</math>0.2</u>
Base + Prior (30M)	24.9 $\pm$ 0.1	30.0 $\pm$ 0.1	<u>21.0<math>\pm</math>0.4</u>	<u>29.8<math>\pm</math>0.3</u>	<u>18.6<math>\pm</math>0.0</u>	19.5 $\pm$ 0.2	–	–

Table: BLEU scores of each NMT method (mean and std of 3 runs)

The hyperparameters of each method were tuned on the DE→EN dev-set:  
LS( $\alpha = 0.1$ ), shallow-fusion( $\beta=0.1$ ), LM-prior ( $\lambda=0.5$ ,  $\tau=2$ )

# Results: Extremely Low-Resource Conditions



**Figure:** BLEU scores on the DE→EN test set with different scales of parallel data. Mean of 3 runs reported.

# Perplexity Scores of LMs

<b>language</b>	<b>3M (PPL↓)</b>	<b>30M (PPL↓)</b>
English	29.70	<b>25.02</b>
German	22.71	<b>19.22</b>
Turkish	<b>22.78</b>	–

**Table:** Perplexity scores for LMs trained on each language's monolingual data, computed on a small held-out validation set per language.

# Connection to Knowledge Distillation

Distillation also uses as **soft targets** the distributions of a teacher model

- Knowledge distillation (word-level) for NMT(Kim and Rush, 2016)

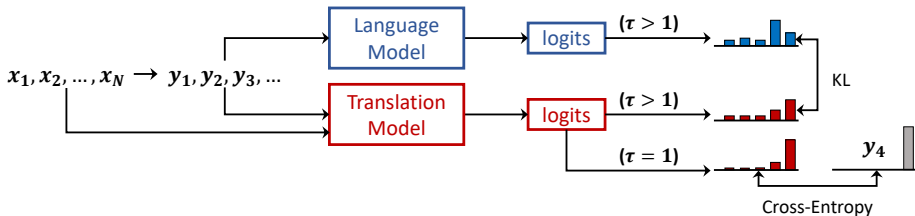
# Connection to Knowledge Distillation

Distillation also uses as **soft targets** the distributions of a teacher model

- Knowledge distillation (word-level) for NMT (Kim and Rush, 2016)

Use softmax-temperature  $\tau$  to reveal **extra information** from LM

- popularized as “*dark knowledge*” by (Hinton et al., 2015)



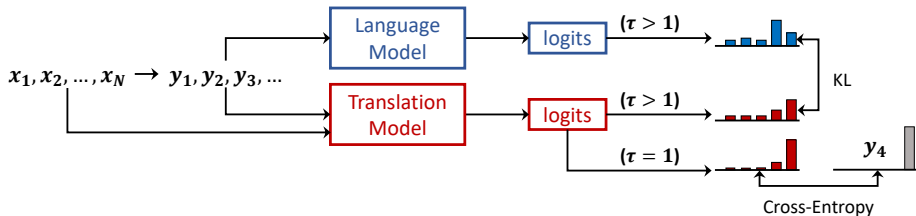
# Connection to Knowledge Distillation

Distillation also uses as **soft targets** the distributions of a teacher model

- Knowledge distillation (word-level) for NMT (Kim and Rush, 2016)

Use softmax-temperature  $\tau$  to reveal **extra information** from LM

- popularized as “*dark knowledge*” by (Hinton et al., 2015)



$$\mathcal{L} = \sum_{t=1}^N -\log p_{\text{TM}}(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{x}) + \lambda \tau^2 D_{\text{KL}}(p_{\text{LM}}(\mathbf{y}_t | \mathbf{y}_{<t}; \tau) \| p_{\text{TM}}(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{x}; \tau))$$

# $\mathcal{L}_{\text{KL}}$ Sensitivity Analysis

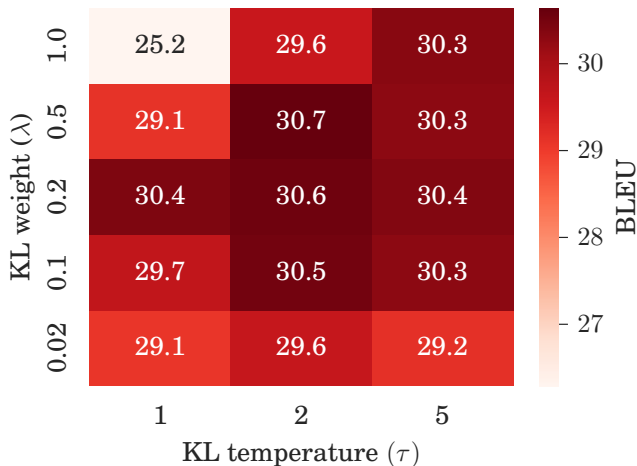
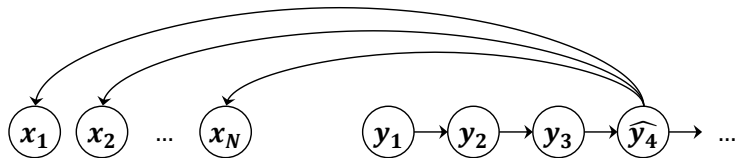


Figure: BLEU scores on the DE $\rightarrow$ EN test set of models trained with different  $\lambda$  and  $\tau$  for the  $\mathcal{L}_{\text{KL}}$ . Mean of 3 runs for each combination reported.

## Statistical Machine Translation

It models the “reverse translation probability”  $p(\mathbf{y}|\mathbf{x}) \propto p(\mathbf{x}|\mathbf{y}) \times p(\mathbf{y})$

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y}} \underbrace{p_{\text{TM}}(\mathbf{x}|\mathbf{y})}_{\text{translation model}} \times \underbrace{p_{\text{LM}}(\mathbf{y})}_{\text{language model}}$$

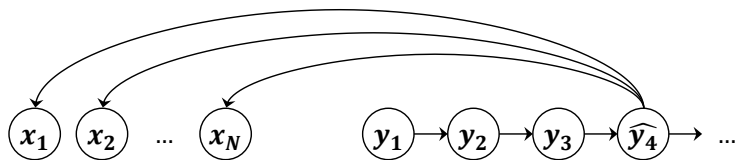




## Statistical Machine Translation

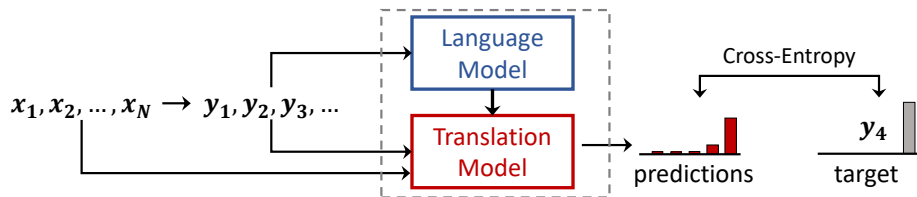
It models the “reverse translation probability”  $p(\mathbf{y}|\mathbf{x}) \propto p(\mathbf{x}|\mathbf{y}) \times p(\mathbf{y})$

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y}} \underbrace{p_{\text{TM}}(\mathbf{x}|\mathbf{y})}_{\text{translation model}} \times \underbrace{p_{\text{LM}}(\mathbf{y})}_{\text{language model}}$$

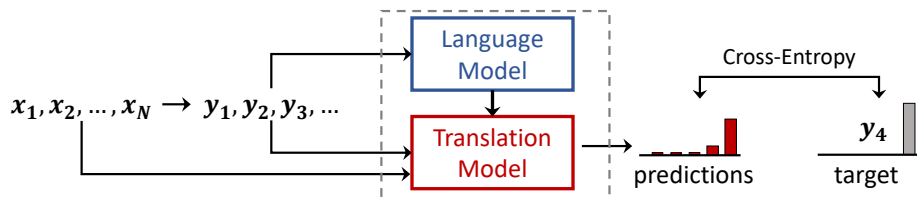


Decoding with seq2seq is **slow!**

# Prior Work: Language Model Fusion



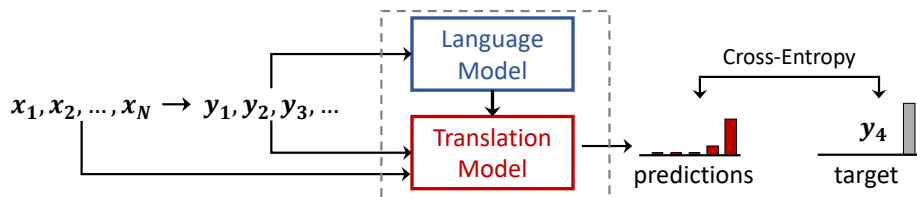
# Prior Work: Language Model Fusion



**Feature-based:** incorporate the representations of the LM

- deep-fusion (Gulcehre et al., 2015)
- cold-fusion (Sriram et al., 2018)

# Prior Work: Language Model Fusion



**Feature-based:** incorporate the representations of the LM

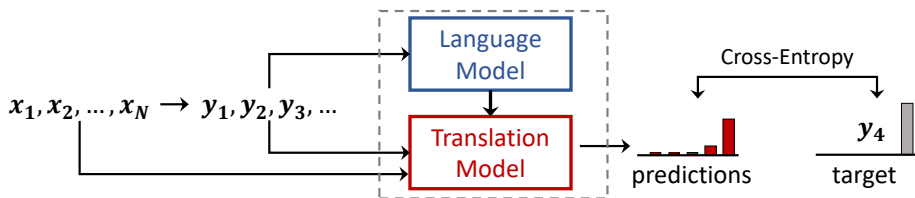
- deep-fusion (Gulcehre et al., 2015)
- cold-fusion (Sriram et al., 2018)

**Probability Interpolation:** incorporate the predictions of the LM

- shallow-fusion (Gulcehre et al., 2015)

$$\hat{y} = \arg \max_{y} \log p_{\text{TM}}(y_t | y_{<t}, x) + \beta \log p_{\text{LM}}(y_t | y_{<t})$$

# Prior Work: Language Model Fusion



**Feature-based:** incorporate the representations of the LM

- deep-fusion (Gulcehre et al., 2015)
- cold-fusion (Sriram et al., 2018)

**Probability Interpolation:** incorporate the predictions of the LM

- shallow-fusion (Gulcehre et al., 2015)

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y}} \log p_{\text{TM}}(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{x}) + \beta \log p_{\text{LM}}(\mathbf{y}_t | \mathbf{y}_{<t})$$

- simple-fusion (Stahlberg et al., 2018)  $\rightarrow$  POSTNORM:

$$p(\mathbf{y}_t) = \text{softmax}(\log p_{\text{TM}}(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{x}) + \log p_{\text{LM}}(\mathbf{y}_t | \mathbf{y}_{<t}))$$

- Caglar Gulcehre, Orhan Firat, Kelvin Xu, Kyunghyun Cho, Loic Barrault, Hwei-Chi Lin, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2015. On using monolingual corpora in neural machine translation. *arXiv preprint arXiv:1503.03535*.
- Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean. 2015. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, abs/1503.02531.
- Yoon Kim and Alexander M. Rush. 2016. Sequence-level knowledge distillation. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 1317–1327, Austin, Texas, USA.
- Anuroop Sriram, Heewoo Jun, Sanjeev Satheesh, and Adam Coates. 2018. Cold fusion: Training seq2seq models together with language models. In *Proceedings of Interspeech*, pages 387–391.
- Felix Stahlberg, James Cross, and Veselin Stoyanov. 2018. Simple fusion: Return of the language model. In *Proceedings of the Conference on Machine Translation*, pages 204–211, Belgium, Brussels.